

# **Bayesian Learning for Robot Audition**

Keynote, HSCMA 2017, San Francisco

#### **Christine Evers**

Imperial College London Dept. Electrical and Electronic Engineering



# Introduction

C. Evers  $\cdot$  Bayesian Learning for Robot Audition  $\cdot$  HSCMA 2017 - 2/45

## Listening in Realistic Environments

Imperial College London



## Listening in Realistic Environments

Imperial College London



# Listening in Realistic Environments

Imperial College London



### Acoustic Scene Mapping

- Estimate the positions of surrounding sources
- Focus and interact with users
- Situational awareness to react to stimuli in the environment



Socially assistive robots. Photo credit: SoftBank Robotics



Search-and-rescue robots. Photo credit: U.S. Army Space and Missile Defense Command

Impact



Smart homes. Photo credit: theappsolutions.com



Hearing aids. Photo credit: Oregon Lions Sight & Hearing Foundation

### Impact



Virtual reality. Photo credit: Samsung

- Aim: Estimate source direction(s) at one time instance
- Generalized cross correlation (GCC), Multiple Signal Classification (MUSIC), Intensity vectors, ...

<sup>&</sup>lt;sup>1</sup>J. H. DiBiase, H. F. Silverman, and M. S. Brandstein, "Robust Localization in Reverberant Rooms," in *Microphone Arrays*, Springer, 2001.

<sup>&</sup>lt;sup>2</sup>J. Benesty, J. Chen, and Y. Huang, *Microphone Array Signal Processing*, Springer 2008.

# Sound Source Localization

- Aim: Estimate source direction(s) at one time instance
- Generalized cross correlation (GCC), Multiple Signal Classification (MUSIC), Intensity vectors, ...

## Localization Challenges



# Localization Challenges



# Localization Challenges



# Localization & Tracking

### Sound Source Localization

• Aim: Estimate source direction(s) at one time instance

### Sound Source Localization

• Aim: Estimate source direction(s) at one time instance

### Sound Source Tracking

- Utilize localization estimates as "measurements"
- Exploit temporal models of the source dynamics
- Estimate smoothed source trajectories over time









# Acoustic Scene Mapping Challenges

Imperial College London



# **Bayesian Inference**

C. Evers  $\cdot$  Bayesian Learning for Robot Audition  $\cdot$  HSCMA 2017 -  $10\,/\,45$ 

# Real-time Scene Mapping

- Sequential estimation required
- Tracking domain:
  - Each window of speech data considered one time step, t
  - Localization estimates at t: Measurements, z<sub>t</sub>
  - Aim: Estimate source position, s<sub>t</sub> from trajectory z<sub>1:t</sub>
- **Probabilistic perspective:** Propagation of posterior density function (pdf)

# Real-time Scene Mapping

- Tracking domain:
  - Each window of speech data considered one time step, t
  - Localization estimates at t: Measurements, z<sub>t</sub>
  - Aim: Estimate source position, s<sub>t</sub> from trajectory z<sub>1:t</sub>
- **Probabilistic perspective:** Propagation of posterior density function (pdf)

$$p(\mathbf{s}_{t-1} \mid \mathbf{z}_{1:t-1}) \xrightarrow{f(\cdot)} p(\mathbf{s}_t \mid \mathbf{z}_{1:t})$$

# Real-time Scene Mapping

- Sequential estimation required
- Tracking domain:
  - Each window of speech data considered one time step, t
  - Localization estimates at t: Measurements, z<sub>t</sub>
  - Aim: Estimate source position, s<sub>t</sub> from trajectory z<sub>1:t</sub>
- **Probabilistic perspective:** Propagation of posterior density function (pdf)

$$p(\mathbf{s}_{t-1} | \mathbf{z}_{1:t-1}) \xrightarrow{f(\cdot)} p(\mathbf{s}_t | \mathbf{z}_{1:t}) \xrightarrow{f(\cdot)} p(\mathbf{s}_t | \mathbf{z}_{1:t}) \xrightarrow{f(\cdot)} p(\mathbf{s}_{t+1} | \mathbf{z}_{1:t+1})$$

Imperial College London

# Bayesian Recipe

### Step 1 - Probability transformation

Dynamical model:	$\mathbf{s}_t = f(\mathbf{s}_{t-1}, \mathbf{v}_t)$	$\stackrel{\text{Prior}}{\Rightarrow}$	$p\left(\left.\mathbf{s}_{t}\right \left.\mathbf{s}_{t-1}\right)\right.$
Measurement model:	$\mathbf{z}_t = h(\mathbf{s}_t, \mathbf{w}_t)$	$\stackrel{\rm Likelihood}{\Rightarrow}$	$p\left(\left.\mathbf{z}_{t}\right \right.\mathbf{s}_{t}\right)$

# Bayesian Recipe

#### Step 1 - Probability transformation

#### Step 2 - Prediction using the prior

$$p\left(\left.\mathbf{s}_{t}\right|\left.\mathbf{z}_{1:t-1}\right)=\int p\left(\left.\mathbf{s}_{t-1}\right|\left.\mathbf{z}_{1:t-1}\right)p\left(\left.\mathbf{s}_{t}\right|\left.\mathbf{s}_{t-1}\right)d\mathbf{s}_{t-1}\right.\right)$$

#### Step 1 - Probability transformation

#### Step 2 - Prediction using the prior

$$p\left(\left.\mathbf{s}_{t}\right|\left.\mathbf{z}_{1:t-1}\right)=\int p\left(\left.\mathbf{s}_{t-1}\right|\left.\mathbf{z}_{1:t-1}\right)p\left(\left.\mathbf{s}_{t}\right|\left.\mathbf{s}_{t-1}\right)d\mathbf{s}_{t-1}\right.\right.$$

#### Step 3 - Bayes's theorem

$$p\left(\left.\mathbf{s}_{t}\right|\left.\mathbf{z}_{1:t}\right)=\frac{p\left(\left.\mathbf{z}_{t}\right|\left.\mathbf{s}_{t},\mathbf{z}_{1:t-1}\right)p\left(\left.\mathbf{s}_{t}\right|\left.\mathbf{z}_{1:t-1}\right)\right)}{p\left(\left.\mathbf{z}_{t}\right|\left.\mathbf{z}_{1:t-1}\right)\right)}$$

# Source DoA Tracking

C. Evers  $\cdot$  Bayesian Learning for Robot Audition  $\cdot$  HSCMA 2017 - 13 / 45

# Tracking and Localization Errors



Dynamical Model of DoA 
$$\boldsymbol{\omega}_{t} \triangleq \begin{bmatrix} \phi_{t}, \theta_{t}, \dot{\phi}_{t}, \dot{\theta}_{t} \end{bmatrix}^{T}$$
  
Source state:  $\boldsymbol{\omega}_{t} = \vartheta \left( \mathbf{F}_{t} \boldsymbol{\omega}_{t-1} + \mathbf{v}_{t} \right)$  (1)  
Constant-velocity dynamics:  $\mathbf{F}_{t} \triangleq \begin{bmatrix} \mathbf{I}_{2} & \Delta_{t} \mathbf{I}_{2} \\ \mathbf{0}_{2 \times 2} & \mathbf{I}_{2} \end{bmatrix}$  (2)  
where  $\vartheta \left( \boldsymbol{\omega}_{t} \right) = \mod \left( \boldsymbol{\omega}_{t}, \left[ 2\pi, \pi, 2\pi, \pi \right]^{T} \right)$  and  $\mathbf{v}_{t}$ : process noise with covariance  $\mathbf{Q}$ .

(3)

Dynamical Model of DoA 
$$\boldsymbol{\omega}_t \triangleq \begin{bmatrix} \phi_t, \theta_t, \dot{\phi}_t, \dot{\theta}_t \end{bmatrix}^T$$
  
Source state:  $\boldsymbol{\omega}_t = \vartheta \left( \mathbf{F}_t \boldsymbol{\omega}_{t-1} + \mathbf{v}_t \right)$  (1)  
Constant-velocity dynamics:  $\mathbf{F}_t \triangleq \begin{bmatrix} \mathbf{I}_2 & \Delta_t \mathbf{I}_2 \\ \mathbf{0}_{2 \times 2} & \mathbf{I}_2 \end{bmatrix}$  (2)  
where  $\vartheta \left( \boldsymbol{\omega}_t \right) = \mod \left( \boldsymbol{\omega}_t, [2\pi, \pi, 2\pi, \pi]^T \right)$  and  $\mathbf{v}_t$ : process noise with  
covariance  $\mathbf{Q}$ .  
Measurement Model of localized DoA  $\mathbf{z}_t \triangleq \begin{bmatrix} \hat{\phi}_t, \hat{\theta}_t \end{bmatrix}^T$ 

Source direction:  $\mathbf{z}_t = \vartheta \left( \mathbf{H} \, \boldsymbol{\omega}_t + \mathbf{w}_t \right)$ 

where 
$$\mathbf{H} \triangleq \begin{bmatrix} \mathbf{I}_2 & \mathbf{0}_{2 \times 2} \\ \mathbf{0}_{2 \times 4} \end{bmatrix}$$
 and  $\mathbf{w}_t$ : measurement noise with covariance  $\mathbf{R}$ .

 $|\phi_t, \theta_t|$ 

# Single Source Tracking - Kalman Filter

For source constrained to area where  $2\pi$  crossing is avoided:

- Approximate by linear state space
- Prior and likelihood: Gaussian

# Single Source Tracking - Kalman Filter

For source constrained to area where  $2\pi$  crossing is avoided:

- Approximate by linear state space
- Prior and likelihood: Gaussian

Prediction - Marginalization with Gaussian prior:

$$p\left(\boldsymbol{\omega}_{t} \mid \mathbf{z}_{1:t-1}\right) = \mathcal{N}\left(\boldsymbol{\omega}_{t} \mid \mathbf{m}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}\right)$$

where  $\mathbf{m}_{t|t-1} = \mathbf{F}_t \, \mathbf{m}_{t-1}$  and  $\boldsymbol{\Sigma}_{t|t-1} = \mathbf{F}_t \, \boldsymbol{\Sigma}_{t-1} \, \mathbf{F}_t^T + \mathbf{Q}_t$ 

# Single Source Tracking - Kalman Filter

For source constrained to area where  $2\pi$  crossing is avoided:

- Approximate by linear state space
- Prior and likelihood: Gaussian

Prediction - Marginalization with Gaussian prior:

$$p\left(\boldsymbol{\omega}_{t} \mid \mathbf{z}_{1:t-1}\right) = \mathcal{N}\left(\boldsymbol{\omega}_{t} \mid \mathbf{m}_{t|t-1}, \boldsymbol{\Sigma}_{t|t-1}\right)$$

where 
$$\mathbf{m}_{t|t-1} = \mathbf{F}_t \, \mathbf{m}_{t-1}$$
 and  $\boldsymbol{\Sigma}_{t|t-1} = \mathbf{F}_t \, \boldsymbol{\Sigma}_{t-1} \, \mathbf{F}_t^T + \mathbf{Q}_t$ 

Update - Bayes's theorem for Gaussian likelihood:

$$p\left(\boldsymbol{\omega}_{t} \mid \mathbf{z}_{1:t}\right) = \mathcal{N}\left(\boldsymbol{\omega}_{t} \mid \mathbf{m}_{t}, \boldsymbol{\Sigma}_{t}\right)$$

where  $\mathbf{m}_t = \mathbf{m}_{t|t-1} + \mathbf{K}_t \left( \mathbf{z}_t - \mathbf{H} \mathbf{m}_{t|t-1} \right)$  and  $\mathbf{\Sigma}_t = \left( \mathbf{I}_4 - \mathbf{K}_t \mathbf{H} \right) \mathbf{\Sigma}_{t|t-1}$ .










# Classical DoA Tracking - Kalman Filter++

- Missing detections due to speech inactivity
- Sources may enter / leave the scene
- Initial source position unknown a priori



# **Multi-Source Tracking**

### Tracking, spurious & missing detections



### Multi-Source Localization



### Multi-Source Localization



- Time-varying number of unlabelled measurements
- Unknown, time-varying number of unknown source states

### Multi-Hypothesis Tracking (MHT)

- **Hypothesis tree:** Enumerate exhaustively all association hypotheses across time
- Final time step: Backward-propagation to identify most likely path
- Computationally prohibitive for real-time processing

<sup>&</sup>lt;sup>1</sup>S. S. Blackman, "Multiple hypothesis tracking for multiple target tracking," IEEE Aerosp. Electr. Sys. Mag., 2004.

<sup>&</sup>lt;sup>2</sup>Y. Bar-Shalom, F. Daum, and J. Huang, "The probabilistic data association filter," IEEE Control Syst. Mag., 2010.

### Multi-Hypothesis Tracking (MHT)

- **Hypothesis tree:** Enumerate exhaustively all association hypotheses across time
- Final time step: Backward-propagation to identify most likely path
- Computationally prohibitive for real-time processing

### Joint Probabilistic Data Association (JPDA)

- At each time step consider all association hypotheses, including spurious detections
- Use pruning after each time step to reduce computational complexity

<sup>&</sup>lt;sup>1</sup>S. S. Blackman, "Multiple hypothesis tracking for multiple target tracking," IEEE Aerosp. Electr. Sys. Mag., 2004.

<sup>&</sup>lt;sup>2</sup>Y. Bar-Shalom, F. Daum, and J. Huang, "The probabilistic data association filter," IEEE Control Syst. Mag., 2010.

### Multi-Hypothesis Tracking (MHT)

- **Hypothesis tree:** Enumerate exhaustively all association hypotheses across time
- Final time step: Backward-propagation to identify most likely path
- Computationally prohibitive for real-time processing

### Joint Probabilistic Data Association (JPDA)

- At each time step consider all association hypotheses, including spurious detections
- Use pruning after each time step to reduce computational complexity
- Estimation of the number of sources typically selected heuristically

<sup>&</sup>lt;sup>1</sup>S. S. Blackman, "Multiple hypothesis tracking for multiple target tracking," IEEE Aerosp. Electr. Sys. Mag., 2004.

<sup>&</sup>lt;sup>2</sup>Y. Bar-Shalom, F. Daum, and J. Huang, "The probabilistic data association filter," IEEE Control Syst. Mag., 2010.

### Multi-Source Tracking - Set Model

### Random Finite Set (RFS) of sources

$$\boldsymbol{\Omega}_{t} = \left[\bigcup_{n=1}^{N_{t-1}} P(\boldsymbol{\omega}_{t-1,n})\right] \cup B_{t}.$$
(4)

$$P(\boldsymbol{\omega}_{t-1,n}) = \begin{cases} \boldsymbol{\omega}_{t,n}, & \text{if } \boldsymbol{\omega}_{t-1,n} \text{ persists} \\ \boldsymbol{\emptyset}, & \text{otherwise.} \end{cases}$$
(5)

where  $|\mathbf{\Omega}_t| = N_t$  and  $B_t$  models the birth of previously inactive sources.

### Multi-Source Tracking - Set Model

Imperial College London

### Random Finite Set (RFS) of sources

$$\boldsymbol{\Omega}_{t} = \begin{bmatrix} \sum_{n=1}^{N_{t-1}} P(\boldsymbol{\omega}_{t-1,n}) \end{bmatrix} \cup B_{t}.$$
(4)

$$P(\boldsymbol{\omega}_{t-1,n}) = \begin{cases} \boldsymbol{\omega}_{t,n}, & \text{if } \boldsymbol{\omega}_{t-1,n} \text{ persists} \\ \boldsymbol{\emptyset}, & \text{otherwise.} \end{cases}$$
(5)

where  $|\mathbf{\Omega}_t| = N_t$  and  $B_t$  models the birth of previously inactive sources.

### Detection RFS of sources and early reflections

$$\mathbf{Z}_{t} = \left[\bigcup_{n=1}^{N_{t}} D(\boldsymbol{\omega}_{t,n})\right] \cup C_{t},$$
(6)

$$D(\boldsymbol{\omega}_{t,n}) = \begin{cases} \mathbf{z}_{t,m}, & \text{if } \boldsymbol{\omega}_{t,n} \text{ detected} \\ \emptyset, & \text{otherwise} \end{cases}$$
(7)

where  $|\mathbf{Z}_t| = M_t$  and  $C_t$  models spurious detections.

Imperial College London

• Number of sources,  $N_t$ , in RFS  $\boldsymbol{\Omega}_t$  is time-varying and unknown

$$\int p(\mathbf{\Omega}_t) \, \delta \mathbf{\Omega} \stackrel{?}{=}$$

<sup>&</sup>lt;sup>1</sup>R. P. S. Mahler, "Multitarget Bayes filtering via first-order multitarget moments," IEEE Tran. Aerosp. Electr. Sys, 2003.

Imperial College London

• Number of sources,  $N_t$ , in RFS  $\boldsymbol{\Omega}_t$  is time-varying and unknown

$$\int p(\boldsymbol{\Omega}_t) \, \delta \boldsymbol{\Omega} \stackrel{?}{=} p(\boldsymbol{\emptyset}).$$

<sup>&</sup>lt;sup>1</sup>R. P. S. Mahler, "Multitarget Bayes filtering via first-order multitarget moments," IEEE Tran. Aerosp. Electr. Sys, 2003.

Imperial College London

• Number of sources,  $N_t$ , in RFS  $\boldsymbol{\Omega}_t$  is time-varying and unknown

$$\int p(\boldsymbol{\Omega}_t) \, \delta \boldsymbol{\Omega} \stackrel{?}{=} p(\boldsymbol{\emptyset}) + \int p(\left\{\boldsymbol{\omega}_{t,1}\right\}) d\boldsymbol{\omega}_{t,1}.$$

<sup>&</sup>lt;sup>1</sup>R. P. S. Mahler, "Multitarget Bayes filtering via first-order multitarget moments," IEEE Tran. Aerosp. Electr. Sys, 2003.

Imperial College London

• Number of sources,  $N_t$ , in RFS  $\boldsymbol{\Omega}_t$  is time-varying and unknown

$$\int p(\boldsymbol{\Omega}_t) \, \delta \boldsymbol{\Omega} = p(\boldsymbol{\emptyset}) + \sum_{n=1}^\infty \frac{1}{n!} \int \ldots \int p(\left\{\boldsymbol{\omega}_{t,1}, \ldots, \boldsymbol{\omega}_{t,n}\right\}) d\boldsymbol{\omega}_{t,1} \ldots d\boldsymbol{\omega}_{t,n}.$$

<sup>&</sup>lt;sup>1</sup>R. P. S. Mahler, "Multitarget Bayes filtering via first-order multitarget moments," IEEE Tran. Aerosp. Electr. Sys, 2003.

Imperial College London

• Number of sources,  $N_t$  , in RFS  $\boldsymbol{\Omega}_t$  is time-varying and unknown

$$\int p(\boldsymbol{\Omega}_t)\,\delta\boldsymbol{\Omega} = p(\boldsymbol{\emptyset}) + \sum_{n=1}^\infty \frac{1}{n!}\int \dots \int p(\left\{\boldsymbol{\omega}_{t,1},\dots,\boldsymbol{\omega}_{t,n}\right\})d\boldsymbol{\omega}_{t,1}\dots d\boldsymbol{\omega}_{t,n}.$$

- Multi-source pdf is combinatorially intractable
- Approximate multi-source pdf by its first order moment:
  - Probability hypothesis density (PHD),  $\lambda(\boldsymbol{\omega}_t \mid \mathbf{Z}_{1:t})$
  - Probability that one of the sources has state ω<sub>t</sub>

<sup>&</sup>lt;sup>1</sup>R. P. S. Mahler, "Multitarget Bayes filtering via first-order multitarget moments," IEEE Tran. Aerosp. Electr. Sys, 2003.

Imperial College London

• Number of sources,  $N_t$  , in RFS  $\boldsymbol{\Omega}_t$  is time-varying and unknown

$$\int p(\boldsymbol{\Omega}_t) \, \delta \boldsymbol{\Omega} = p(\boldsymbol{\emptyset}) + \sum_{n=1}^\infty \frac{1}{n!} \int \ldots \int p(\left\{\boldsymbol{\omega}_{t,1}, \ldots, \boldsymbol{\omega}_{t,n}\right\}) d\boldsymbol{\omega}_{t,1} \ldots d\boldsymbol{\omega}_{t,n}.$$

- Multi-source pdf is combinatorially intractable
- Approximate multi-source pdf by its first order moment:
  - Probability hypothesis density (PHD),  $\lambda(\boldsymbol{\omega}_t | \mathbf{Z}_{1:t})$
  - Probability that one of the sources has state ω<sub>t</sub>

$$\int \lambda\left(\left.\boldsymbol{\omega}_{t}\right|\left.\mathbf{Z}_{1:t}\right)d\boldsymbol{\omega}_{t}=\mathbb{E}\left[N_{t}\right]$$

<sup>&</sup>lt;sup>1</sup>R. P. S. Mahler, "Multitarget Bayes filtering via first-order multitarget moments," IEEE Tran. Aerosp. Electr. Sys, 2003.

Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_t\right) + \lambda^{\mathsf{missed}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t-1}\right) + \lambda^{\mathsf{detect}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t}\right)\right.$$

<sup>&</sup>lt;sup>1</sup>C. Evers *et al.*, "Bearing-only Acoustic Tracking of Moving Speakers for Robot Audition", Proc. IEEE DSP, Singapore, 2015

Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_t \right) + \lambda^{\mathsf{missed}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_{1:t-1} \right) + \lambda^{\mathsf{detect}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_{1:t} \right)$$



Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_t\right) + \lambda^{\mathsf{missed}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t-1}\right) + \lambda^{\mathsf{detect}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t}\right)\right.$$



Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_t \right| \left.\mathbf{Z}_t\right) + \lambda^{\mathsf{missed}}\left(\left.\mathbf{s}_t \right| \left.\mathbf{Z}_{1:t-1}\right) + \lambda^{\mathsf{detect}}\left(\left.\mathbf{s}_t \right| \left.\mathbf{Z}_{1:t}\right)\right.$$

Newborn sources: Each DoA may originate from a new source

$$\lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_{t}\right|\left.\mathbf{Z}_{t}\right)=\sum_{m=1}^{M_{t}}p_{b}\,\mathcal{N}\left(\mathbf{s}_{t}\left|\mathbf{z}_{t,m},\,\mathbf{R}\right.\right)$$

where  $p_b$  is the probability of source birth.

<sup>&</sup>lt;sup>1</sup>C. Evers et al., "Bearing-only Acoustic Tracking of Moving Speakers for Robot Audition", Proc. IEEE DSP, Singapore, 2015

Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_t \right) + \lambda^{\mathsf{missed}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_{1:t-1} \right) + \lambda^{\mathsf{detect}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_{1:t} \right)$$



Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_t \right) + \lambda^{\mathsf{missed}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_{1:t-1} \right) + \lambda^{\mathsf{detect}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_{1:t} \right)$$



Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_t\right) + \lambda^{\mathsf{missed}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t-1}\right) + \lambda^{\mathsf{detect}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t}\right)\right.$$

Missed detections: Any previously detected source may be missed:

$$\lambda^{\text{missed}}\left(\left.\mathbf{s}_{t}\right|\left.\mathbf{Z}_{1:t-1}\right)=\sum_{j=1}^{J_{t|t-1}}w_{t|t-1}^{(j)}\mathcal{N}\left(\mathbf{s}_{t}\left|\mathbf{m}_{t|t-1}^{(j)},\mathbf{\Sigma}_{t|t-1}^{(j)}\right.\right)$$

where  $\mathbf{m}_{t|t-1}$  and  $\mathbf{\Sigma}_{t|t-1}$ : Predicted Kalman filter mean / covariance

<sup>&</sup>lt;sup>1</sup>C. Evers et al., "Bearing-only Acoustic Tracking of Moving Speakers for Robot Audition", Proc. IEEE DSP, Singapore, 2015

Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_t\right) + \lambda^{\mathsf{missed}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t-1}\right) + \lambda^{\mathsf{detect}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t}\right)\right.$$

Missed detections: Any previously detected source may be missed:

$$\begin{split} \lambda^{\text{missed}} \left( \left. \mathbf{s}_{t} \right| \, \mathbf{Z}_{1:t-1} \right) &= \sum_{j=1}^{J_{t|t-1}} w_{t|t-1}^{(j)} \, \mathcal{N} \left( \mathbf{s}_{t} \, \big| \, \mathbf{m}_{t|t-1}^{(j)}, \, \mathbf{\Sigma}_{t|t-1}^{(j)} \right) \\ & w_{t|t-1}^{(j)} = p_{s} \left( 1 - p_{d} \right) w_{t-1}^{(j)} \end{split}$$

where  $\mathbf{m}_{t|t-1}$  and  $\mathbf{\Sigma}_{t|t-1}$ : Predicted Kalman filter mean / covariance

<sup>&</sup>lt;sup>1</sup>C. Evers *et al.*, "Bearing-only Acoustic Tracking of Moving Speakers for Robot Audition", Proc. IEEE DSP, Singapore, 2015

Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_t \right) + \lambda^{\mathsf{missed}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_{1:t-1} \right) + \lambda^{\mathsf{detect}} \left( \left. \mathbf{s}_t \right| \, \mathbf{Z}_{1:t} \right)$$



Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_t\right) + \lambda^{\mathsf{missed}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t-1}\right) + \lambda^{\mathsf{detect}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t}\right)\right.$$



Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_t\right) + \lambda^{\mathsf{missed}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t-1}\right) + \lambda^{\mathsf{detect}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t}\right)\right.$$



Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_t\right) + \lambda^{\mathsf{missed}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t-1}\right) + \lambda^{\mathsf{detect}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t}\right)\right.$$

Detected sources: Any DoA may originate from the source

$$\lambda^{\text{detect}}\left(\left.\mathbf{s}_{t}\right|\left.\mathbf{Z}_{1:t}\right) = \sum_{m=1}^{M_{t}} p_{d} \sum_{j=1}^{J_{t|t-1}} w_{t}^{(j)} \mathcal{N}\left(\mathbf{s}_{t}\left|\mathbf{m}_{t}^{(j,m)}, \mathbf{\Sigma}_{t}^{(j,m)}\right.\right)$$

<sup>&</sup>lt;sup>1</sup>C. Evers et al., "Bearing-only Acoustic Tracking of Moving Speakers for Robot Audition", Proc. IEEE DSP, Singapore, 2015

Explicitly model sources, missing detections, and spurious detections:

$$\lambda(\mathbf{s}_t | \mathbf{Z}_{1:t}) = \lambda^{\mathsf{birth}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_t\right) + \lambda^{\mathsf{missed}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t-1}\right) + \lambda^{\mathsf{detect}}\left(\left.\mathbf{s}_t \mid \mathbf{Z}_{1:t}\right)\right.$$

Detected sources: Any DoA may originate from the source

$$w_t^{(j,m)} = p_d \frac{\mathcal{N}\left(\mathbf{z}_{t,m} \, \big| \mathbf{H}\mathbf{m}_{t|t-1}^{(j)}, \mathbf{S}_t^{(j)}\right)}{\kappa\left(\mathbf{z}_{t,m}\right) + \sum_{i=1}^{J_{t|t-1}} \mathcal{N}\left(\mathbf{z}_{t,m} \, \big| \mathbf{H}\mathbf{m}_{t|t-1}^{(i)}, \mathbf{S}_t^{(j)}\right)}$$

where 
$$\mathbf{S}_{t}^{(j)} \triangleq \mathbf{H} \mathbf{\Sigma}_{t|t-1} \mathbf{H}^{T} + \mathbf{R}$$
 and  $\kappa \left( \mathbf{z}_{t,m} \right) \triangleq \mathsf{FAR} \times \mathcal{U} \left[ \mathsf{FoV} \right]$ 

<sup>&</sup>lt;sup>1</sup>C. Evers et al., "Bearing-only Acoustic Tracking of Moving Speakers for Robot Audition", Proc. IEEE DSP, Singapore, 2015

### GMMs for Audio-Visual Fusion

- Audio: Speech inactivity, reverberation / noise, interference
- Vision: Limited field-of-view, visual obstructions, lighting
- Complementary modalities

#### Multi-Modal Measurements

Visual face detector:

Sound source localization:

$$\begin{split} z_t^{(\upsilon)} &= g_\upsilon \left( \mathbf{x}_t, \mathbf{w}_t^{(\upsilon)} \right) \\ z_t^{(a)} &= g_a \left( \mathbf{x}_t, \mathbf{w}_t^{(a)} \right) \end{split}$$

#### Find out the details...

Friday, March 3, 13:00 - 14:40 (BATGIRL): ROBOT.5: Audio-Visual Tracking by Density Approximation in a Sequential Bayesian Filtering Framework

- Broadband nature of speech signals
- Peak picking difficult
- Each combination of  ${\cal M}_t$  localized DoAs may contain information about the source

### Multi-detection likelihood

$$p\left(\left.\mathbf{Z}_{t}\right|\right.\boldsymbol{\omega}_{t}\right) = \sum_{\mathbf{P} \boxminus \mathbf{Z}_{t}} \prod_{\ell \in \left(\mathbf{Z}_{t} - \mathbf{P}\right)} \kappa\left(\mathbf{z}(t,\ell)\right) \prod_{p \in \mathbf{P}} p\left(\left.\mathbf{z}(t,p)\right|\right.\boldsymbol{\omega}_{t}\right),$$

#### Find out the details...

Friday, March 3, 15:10 - 16:50 (BATGIRL): LOC.4: Speaker Tracking in Reverberant Environments using Multiple Directions of Arrival
# **Source Position Tracking**

C. Evers · Bayesian Learning for Robot Audition · HSCMA 2017 - 28 / 45

## Tracking Source Positions



## Source Position Tracking - Model

Imperial College London

Dynamical Model of 
$$\mathbf{s}_t \triangleq \left[x, y, z, \dot{x}, \dot{y}, \dot{z}\right]^T$$

Source state:  $\mathbf{s}_t = \mathbf{D}_t \mathbf{s}_{t-1} + \mathbf{v}_t$  (8)

where  $\mathbf{D}_t$ : Langevin model<sup>1</sup>, and  $\Delta_t$  is the time step

Measurement Model of 
$$\mathbf{z}_t \triangleq \left[\hat{\phi}_t, \hat{\theta}_t\right]^T$$

Source direction:  $\mathbf{z}_t = h(\mathbf{s}_t)$ 

$$h(\mathbf{s}_t) + \mathbf{w}_t$$
 (9)

where  $h(\cdot)$  is the Cartesian-to-spherical transformation.

<sup>&</sup>lt;sup>1</sup>E. A. Lehmann and R. C. Williamson, "Particle Filter Design Using Importance Sampling for Acoustic Source Localisation and Tracking in Reverberant Environments," EURASIP J. Adv. Signal Process., 2006.

#### Estimate 6D states from 2D measurements

## Unmeasured Range Inference

#### Estimate 6D states from 2D measurements

Kalman filter (KF) update equation

$$\boldsymbol{\mu}_{t} = \boldsymbol{\mu}_{t|t-1} + \mathbf{K}_{t} \left( \mathbf{z}_{t} - h \left( \boldsymbol{\mu}_{t|t-1} \right) \right)$$

where  $\mu_{t|t-1}$  and  $\mu_t$  are the predicted / updated KF mean of  $\mathbf{s}_t$ .

## Unmeasured Range Inference

#### Estimate 6D states from 2D measurements

Kalman filter (KF) update equation

$$\underbrace{\underline{\mu}_t}_{6\times 1} = \underbrace{\underline{\mu}_{t|t-1}}_{6\times 1} + \underbrace{\mathbf{K}_t}_{6\times 2} (\underbrace{\mathbf{z}_t}_{2\times 1} - \underbrace{h(\underline{\mu}_{t|t-1}}_{2\times 1}))$$

where  $\boldsymbol{\mu}_{t|t-1}$  and  $\boldsymbol{\mu}_t$  are the predicted / updated KF mean of  $\mathbf{s}_t.$ 

## Unmeasured Range Inference

#### Estimate 6D states from 2D measurements

Kalman filter (KF) update equation

$$\underbrace{\underline{\mu}_t}_{6\times 1} = \underbrace{\underline{\mu}_{t|t-1}}_{6\times 1} + \underbrace{\mathbf{K}_t}_{6\times 2} (\underbrace{\mathbf{z}_t}_{2\times 1} - \underbrace{h(\underline{\mu}_{t|t-1}}_{2\times 1}))$$

where  $\mu_{t|t-1}$  and  $\mu_t$  are the predicted / updated KF mean of  $\mathbf{s}_t$ .

- Kalman gain: Map between state and measurement space
- Pre-populate KF mean with range hypothesis
- New directional information updated through measurements
- Range component is extrapolated through time

Imperial College London



Imperial College London



Imperial College London



Imperial College London



# **Acoustic SLAM**



## Simultaneous Localization and Mapping



<sup>&</sup>lt;sup>1</sup>C. Evers et al., "Acoustic Simultaneous Localization and Mapping (a-SLAM) of a Moving Microphone Array and its Surrounding Speakers", Proc. ICASSP, Shanghai, 2016.

<sup>&</sup>lt;sup>2</sup>C. Evers *et al.*, "Localization of Moving Microphone Arrays from Moving Sound Sources for Robot Audition", Proc. EUSIPCO, Budapest, 2016.

<sup>&</sup>lt;sup>3</sup>C. Evers, A. H. Moore, and P. A. Naylor, "Towards Informative Path Planning for Acoustic SLAM", Proc. DAGA, Aachen, 2015.

## Simultaneous Localization and Mapping



<sup>&</sup>lt;sup>1</sup>C. Evers et al., "Acoustic Simultaneous Localization and Mapping (a-SLAM) of a Moving Microphone Array and its Surrounding Speakers", Proc. ICASSP, Shanghai, 2016.

<sup>&</sup>lt;sup>2</sup>C. Evers *et al.*, "Localization of Moving Microphone Arrays from Moving Sound Sources for Robot Audition", Proc. EUSIPCO, Budapest, 2016.

<sup>&</sup>lt;sup>3</sup>C. Evers, A. H. Moore, and P. A. Naylor, "Towards Informative Path Planning for Acoustic SLAM", Proc. DAGA, Aachen, 2015.

## Simultaneous Localization and Mapping



- Acoustic Simultaneous Localization & Mapping (aSLAM):
  - Require sensor position to map moving sources
  - Use the scene map to estimate the sensor position / orientation

<sup>&</sup>lt;sup>1</sup>C. Evers *et al.*, "Acoustic Simultaneous Localization and Mapping (a-SLAM) of a Moving Microphone Array and its Surrounding Speakers", Proc. ICASSP, Shanghai, 2016.

<sup>&</sup>lt;sup>2</sup>C. Evers et al., "Localization of Moving Microphone Arrays from Moving Sound Sources for Robot Audition", Proc. EUSIPCO, Budapest, 2016.

<sup>&</sup>lt;sup>3</sup>C. Evers, A. H. Moore, and P. A. Naylor, "Towards Informative Path Planning for Acoustic SLAM", Proc. DAGA, Aachen, 2015.

#### Approach 1: Active SLAM

- Active sensing for room geometry inference and sensor localization
- Cannot estimate source trajectories
- Intrusive / in current form not suitable for human-robot interaction

#### Approach 1: Active SLAM

- Active sensing for room geometry inference and sensor localization
- Cannot estimate source trajectories
- Intrusive / in current form not suitable for human-robot interaction

#### Approach 2: Classic visual-based SLAM

- Fundamental assumption: Static, continuously active sources
- Problematic for speech inactivity
- Prone to divergence for high false alarm rates

- Passive sensing for human-robot interaction
- Exploit directional information gleaned from surrounding sources
- Crucial robustness against reverb, speech inactivity, source motion
- Novel, generalized SLAM framework that jointly estimate:
  - a) A PHD filter for multi-source tracking
  - b) Estimate robot positions, by exploiting source map, robot reports, prior information about the robot motion

<sup>&</sup>lt;sup>1</sup>C. Evers and P. A. Naylor, "Generalized Dynamic Scene Mapping", submitted, IEEE Tran. Sig. Proc.

- Passive sensing for human-robot interaction
- Exploit directional information gleaned from surrounding sources
- Crucial robustness against reverb, speech inactivity, source motion
- Novel, generalized SLAM framework that jointly estimate:
  - a) A PHD filter for multi-source tracking
  - b) Estimate robot positions, by exploiting source map, robot reports, prior information about the robot motion

#### GEneralized Motion SLAM (GEM-SLAM)

$$\lambda\left(\left.\mathbf{S}_{t},\mathbf{r}_{t}\right|\left.\mathbf{Z}_{1:t},\mathbf{y}_{1:t}\right)=\lambda\left(\left.\mathbf{r}_{t}\right|\left.\mathbf{y}_{1:t}\right)\lambda\left(\left.\mathbf{S}_{t}\right|\left.\mathbf{r}_{t},\mathbf{Z}_{1:t}\right)\right.$$

<sup>&</sup>lt;sup>1</sup>C. Evers and P. A. Naylor, "Generalized Dynamic Scene Mapping", submitted, IEEE Tran. Sig. Proc.

- Passive sensing for human-robot interaction
- Exploit directional information gleaned from surrounding sources
- Crucial robustness against reverb, speech inactivity, source motion
- Novel, generalized SLAM framework that jointly estimate:
  - a) A PHD filter for multi-source tracking
  - b) Estimate robot positions, by exploiting source map, robot reports, prior information about the robot motion

#### GEneralized Motion SLAM (GEM-SLAM)

$$\lambda\left(\boldsymbol{\Omega}_{t}, \mathbf{r}_{t} \mid \mathbf{Z}_{1:t}, \mathbf{y}_{1:t}\right) = \lambda\left(\left.\mathbf{r}_{t} \mid \mathbf{y}_{1:t}\right) \underbrace{\lambda\left(\left.\mathbf{S}_{t} \mid \mathbf{r}_{t}, \mathbf{Z}_{1:t}\right)}^{\text{Source PHD}}$$

<sup>&</sup>lt;sup>1</sup>C. Evers and P. A. Naylor, "Generalized Dynamic Scene Mapping", submitted, IEEE Tran. Sig. Proc.

- Passive sensing for human-robot interaction
- Exploit directional information gleaned from surrounding sources
- Crucial robustness against reverb, speech inactivity, source motion
- Novel, generalized SLAM framework that jointly estimate:
  - a) A PHD filter for multi-source tracking
  - b) Estimate robot positions, by exploiting source map, robot reports, prior information about the robot motion

#### GEneralized Motion SLAM (GEM-SLAM)

$$\lambda\left(\left.\boldsymbol{\Omega}_{t}, \mathbf{r}_{t} \right| \left.\mathbf{Z}_{1:t}, \mathbf{y}_{1:t}\right) = \overbrace{\lambda\left(\left.\mathbf{r}_{t}\right| \left.\mathbf{y}_{1:t}\right)}^{\text{Robot PHD}} \lambda\left(\left.\mathbf{S}_{t}\right| \left.\mathbf{r}_{t}, \mathbf{Z}_{1:t}\right)\right.$$

<sup>&</sup>lt;sup>1</sup>C. Evers and P. A. Naylor, "Generalized Dynamic Scene Mapping", submitted, IEEE Tran. Sig. Proc.



Executed robot state, 
$$\mathbf{r}_{t} \triangleq \begin{bmatrix} x_{t}, y_{t}, z_{t}, v_{t}, \gamma_{t} \end{bmatrix}^{T} = \begin{bmatrix} \mathbf{p}_{t}^{T}, \gamma_{t} \end{bmatrix}^{T}$$
:  
 $\mathbf{p}_{t} = \mathbf{G}(\gamma_{t}) \mathbf{p}_{t-1} + \mathbf{v}_{t,\mathbf{p}}$   
 $\gamma_{t} = \vartheta(\gamma_{t-1} + v_{t,\gamma})$ 
 $\mathbf{G}(\gamma_{t}) \triangleq \begin{bmatrix} 1 & 0 & 0 & \Delta_{t} \sin \gamma_{t} \\ 0 & 1 & 0 & \Delta_{t} \cos \gamma_{t} \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$   
where  $\mathbf{p}_{t} \triangleq \begin{bmatrix} x_{t}, y_{t}, z_{t}, v_{t} \end{bmatrix}^{T}$  with speed,  $v_{t}$ , and orientation,  $\gamma_{t}$ .

Reported robot state: 
$$\mathbf{y}_t \triangleq \begin{bmatrix} y_{t,v} & y_{t,\gamma} \end{bmatrix}^T$$
 $y_{t,v} = v_t + w_{t,v},$  $w_{t,v} \sim \mathcal{N}\left(0, \sigma_{w_{t,v}}^2\right)$  $y_{t,\gamma} = \vartheta\left(\gamma_t + w_{t,\gamma}\right),$  $w_{t,\gamma} \sim \mathcal{N}\left(0, \sigma_{w_{t,\gamma}}^2\right)$ 



Reported robot state: $\mathbf{y}_t \triangleq \big[y_{t,v}$	$[y_{t,\gamma}]^T$
$y_{t,v} = v_t + w_{t,v}, $	$w_{t,v} \sim \mathcal{N}\left(0,  \sigma^2_{w_{t,v}}\right)$
$\boldsymbol{y}_{t,\gamma}=\vartheta\left(\gamma_t+\boldsymbol{w}_{t,\gamma}\right),$	$w_{t,\gamma} \sim \mathcal{N}\left(0,  \sigma^2_{w_{t,\gamma}}\right)$

Non-linear orientation and dependency of position on orientation

- Particle filter to approximate robot density
- Classic SLAM: Prior importance sampling

## GEM-SLAM - Optimized Particle Filter

Imperial College London

Imperial College London

Optimized particle filter of 
$$\mathbf{\hat{r}}_{t}^{(i)} = \left[ \left[ \mathbf{\hat{p}}_{t}^{(i)} 
ight]^{T}, \mathbf{\hat{\gamma}}_{t}^{(i)} 
ight]^{T}$$

- Sampling of orientation,  $\hat{\gamma}_t^{(i)}$  , from wrapped Kalman filter:
- For each  $\hat{\gamma}_t^{(i)}$  , optimal sampling of pose,  $\hat{\mathbf{p}}_t^{(i)}$  , from Kalman filter

$$\lambda\left(\left.\mathbf{r}_{t}\right|\left.\mathbf{y}_{1:t},\mathbf{Z}_{t}\right)\approx\sum_{i=1}^{L}\alpha_{t}^{\left(i\right)}\delta_{\hat{\mathbf{r}}_{t}^{\left(i\right)}}\left(\mathbf{r}_{t}\right).$$

Imperial College London

Optimized particle filter of 
$$\mathbf{\hat{r}}_{t}^{(i)} = \left[ \left[ \mathbf{\hat{p}}_{t}^{(i)} 
ight]^{T}, \mathbf{\hat{\gamma}}_{t}^{(i)} 
ight]^{T}$$

- Sampling of orientation,  $\hat{\gamma}_t^{(i)}$  , from wrapped Kalman filter:
- For each  $\hat{\gamma}_t^{(i)}$  , optimal sampling of pose,  $\hat{\mathbf{p}}_t^{(i)}$  , from Kalman filter

$$\lambda\left(\left.\mathbf{r}_{t}\right|\left.\mathbf{y}_{1:t},\mathbf{Z}_{t}\right)\approx\sum_{i=1}^{L}\alpha_{t}^{\left(i\right)}\delta_{\hat{\mathbf{r}}_{t}^{\left(i\right)}}\left(\mathbf{r}_{t}\right).$$

#### Particle weights

$$\alpha_t^{(i)} = \alpha_{t-1}^{(i)} \, p(\mathbf{z}_t | \hat{\gamma}_t^{(i)}, \hat{\mathbf{p}}_t^{(i)}) \, \mathcal{L}(\mathbf{Z}_t | \mathbf{r}_t^{(i)})$$

Optimized particle filter of 
$$\mathbf{\hat{r}}_{t}^{(i)} = \left[ \left[ \mathbf{\hat{p}}_{t}^{(i)} 
ight]^{T}, \mathbf{\hat{\gamma}}_{t}^{(i)} 
ight]^{T}$$

- Sampling of orientation,  $\hat{\gamma}_t^{(i)}$  , from wrapped Kalman filter:
- For each  $\hat{\gamma}_t^{(i)}$  , optimal sampling of pose,  $\hat{\mathbf{p}}_t^{(i)}$  , from Kalman filter

$$\lambda\left(\left.\mathbf{r}_{t}\right|\left.\mathbf{y}_{1:t},\mathbf{Z}_{t}\right)\approx\sum_{i=1}^{L}\alpha_{t}^{\left(i\right)}\delta_{\hat{\mathbf{r}}_{t}^{\left(i\right)}}\left(\mathbf{r}_{t}\right).$$

#### Particle weights

$$\alpha_t^{(i)} = \underbrace{\alpha_{t-1}^{(i)}}_{\text{Previous particle weight}} p(\mathbf{y}_t | \hat{\gamma}_t^{(i,p)}, \hat{\mathbf{p}}_t^{(i,p)}) \, \mathcal{L}(\mathbf{Z}_t | \mathbf{r}_t^{(i,p)})$$

Optimized particle filter of 
$$\mathbf{\hat{r}}_{t}^{(i)} = \left[ \left[ \mathbf{\hat{p}}_{t}^{(i)} 
ight]^{T}, \mathbf{\hat{\gamma}}_{t}^{(i)} 
ight]^{T}$$

- Sampling of orientation,  $\hat{\gamma}_t^{(i)}$  , from wrapped Kalman filter:
- For each  $\hat{\gamma}_t^{(i)}$  , optimal sampling of pose,  $\hat{\mathbf{p}}_t^{(i)}$  , from Kalman filter

$$\lambda\left(\left.\mathbf{r}_{t}\right|\left.\mathbf{y}_{1:t},\mathbf{Z}_{t}\right)\approx\sum_{i=1}^{L}\alpha_{t}^{\left(i\right)}\delta_{\hat{\mathbf{r}}_{t}^{\left(i\right)}}\left(\mathbf{r}_{t}\right).$$

#### Particle weights

$$\boldsymbol{\alpha}_t^{(i)} = \boldsymbol{\alpha}_{t-1}^{(i)} \quad \underline{p}(\mathbf{y}_t | \hat{\boldsymbol{\gamma}}_t^{(i)}, \hat{\mathbf{p}}_t^{(i)}) \quad \mathcal{L}(\mathbf{Z}_t | \mathbf{r}_t^{(i)})$$

Likelihood of robot reports

Optimized particle filter of 
$$\mathbf{\hat{r}}_{t}^{(i)} = \left[ \left[ \mathbf{\hat{p}}_{t}^{(i)} 
ight]^{T}, \widehat{\gamma}_{t}^{(i)} 
ight]^{T}$$

- Sampling of orientation,  $\hat{\gamma}_t^{(i)}$  , from wrapped Kalman filter:
- For each  $\hat{\gamma}_t^{(i)}$  , optimal sampling of pose,  $\hat{\mathbf{p}}_t^{(i)}$  , from Kalman filter

$$\lambda\left(\left.\mathbf{r}_{t}\right|\left.\mathbf{y}_{1:t},\mathbf{Z}_{t}\right)\approx\sum_{i=1}^{L}\alpha_{t}^{\left(i\right)}\delta_{\hat{\mathbf{r}}_{t}^{\left(i\right)}}\left(\mathbf{r}_{t}\right).$$

#### Particle weights

$$\boldsymbol{\alpha}_t^{(i)} = \boldsymbol{\alpha}_{t-1}^{(i)} \, p(\mathbf{y}_t | \hat{\boldsymbol{\gamma}}_t^{(i)}, \hat{\mathbf{p}}_t^{(i)}) \qquad \underbrace{\mathcal{L}(\mathbf{Z}_t | \mathbf{r}_t^{(i)})}_{t}$$

where  $\mathcal{L}\left(\left.\mathbf{Z}_{t} \mid \mathbf{r}_{t}\right) = e^{-N_{t|t-1}-N_{t,c}} \prod_{m=1}^{M_{t}} \left[\kappa\left(\left.\mathbf{z}_{t,m} \mid \mathbf{r}_{t}\right) + p\left(\left.\mathbf{z}_{t,m} \mid \mathbf{r}_{t}, \mathbf{Z}_{1:t-1}\right)\right]\right]$ 

Multi-measurement evidence

Imperial College London

## Probabilistic anchoring by weighting with multi-measurement evidence



Imperial College London

## Probabilistic anchoring by weighting with multi-measurement evidence



## Probabilistic anchoring by weighting with multi-measurement evidence



# **LOCATA** Challenge

C. Evers  $\cdot$  Bayesian Learning for Robot Audition  $\cdot$  HSCMA 2017 - 41 / 45

## LOCATA Challenge & Corpus

Imperial College London









#### IEEE AASP Challenge: LOCalization and TrAcking (LOCATA) H. Löllmann, C. Evers, H. Barfuß, P. A. Naylor, W. Kellermann
Imperial College London









IEEE AASP Challenge: LOCalization and TrAcking (LOCATA) H. Löllmann, C. Evers, H. Barfuß, P. A. Naylor, W. Kellermann

#### Static & Moving Arrays:

- 15-channel linear array
- 12-channel robot head
- 32-channel spherical mh-acoustic eigenmike



Imperial College London









#### IEEE AASP Challenge: LOCalization and TrAcking (LOCATA) H. Löllmann, C. Evers, H. Barfuß, P. A. Naylor, W. Kellermann

#### Static & Moving arrays:

- 15-channel Linear array
- 12-channel robot head
- 32-channel spherical mh-acoustic eigenmike
- 4-channel hearing aids



Imperial College London









IEEE AASP Challenge: LOCalization and TrAcking (LOCATA) H. Löllmann, C. Evers, H. Barfuß, P. A. Naylor, W. Kellermann

#### Static & Moving arrays:

- 15-channel Linear array
- 12-channel robot head
- 32-channel spherical mh-acoustic eigenmike
- 4-channel hearing aids

#### Static loudspeakers



Imperial College London









IEEE AASP Challenge: LOCalization and TrAcking (LOCATA) H. Löllmann, C. Evers, H. Barfuß, P. A. Naylor, W. Kellermann

#### Static & Moving arrays:

- 15-channel Linear array
- 12-channel robot head
- 32-channel spherical mh-acoustic eigenmike
- 2-channel binaural hearing aids

#### Static loudspeakers Moving talkers



Imperial College London









#### IEEE AASP Challenge: LOCalization and TrAcking (LOCATA) H. Löllmann, C. Evers, H. Barfuß, P. A. Naylor, W. Kellermann



C. Evers · Bayesian Learning for Robot Audition · HSCMA 2017 - 43/45











Imperial College London





## Acknowledgement

• The HSCMA Organizing Committee

### • Collaborators and Co-Authors:

- Patrick A. Naylor (Imperial College London)
- Sharon Gannot, Yuval Dorfan (Bar-Ilan University)
- Boaz Rafaely (Ben-Gurion University)
- Radu Horaud, Israel D. Gebru (INRIA)
- Walter Kellermann, Heinrich Löllmann, Hendrik Barfuss, Alexander Schmidt (Friedrich Alexander University)
- Verena Hafner, Heinrich Mellmann, Claas-Norman Ritter (Humboldt University Berlin)
- James R. Hopgood (University of Edinburgh)
- Daniel Clark (Heriot Watt University)

## • Funding

- EPSRC Fellowship "Acoustic Signal Processing and Scene Analysis for Socially Assistive Robots"
- FP7 "Embodied Audition for RobotS" (EARS)