

Learning Dynamic Stream Weights for Multimodal Dynamical System Models

68. Sitzung der ITG-Fachgruppe "Signalverarbeitung und maschinelles Lernen"

Christopher Schymura and Dorothea Kolossa

September 27th, 2019

Audiovisual speaker tracking



Audiovisual speaker tracking

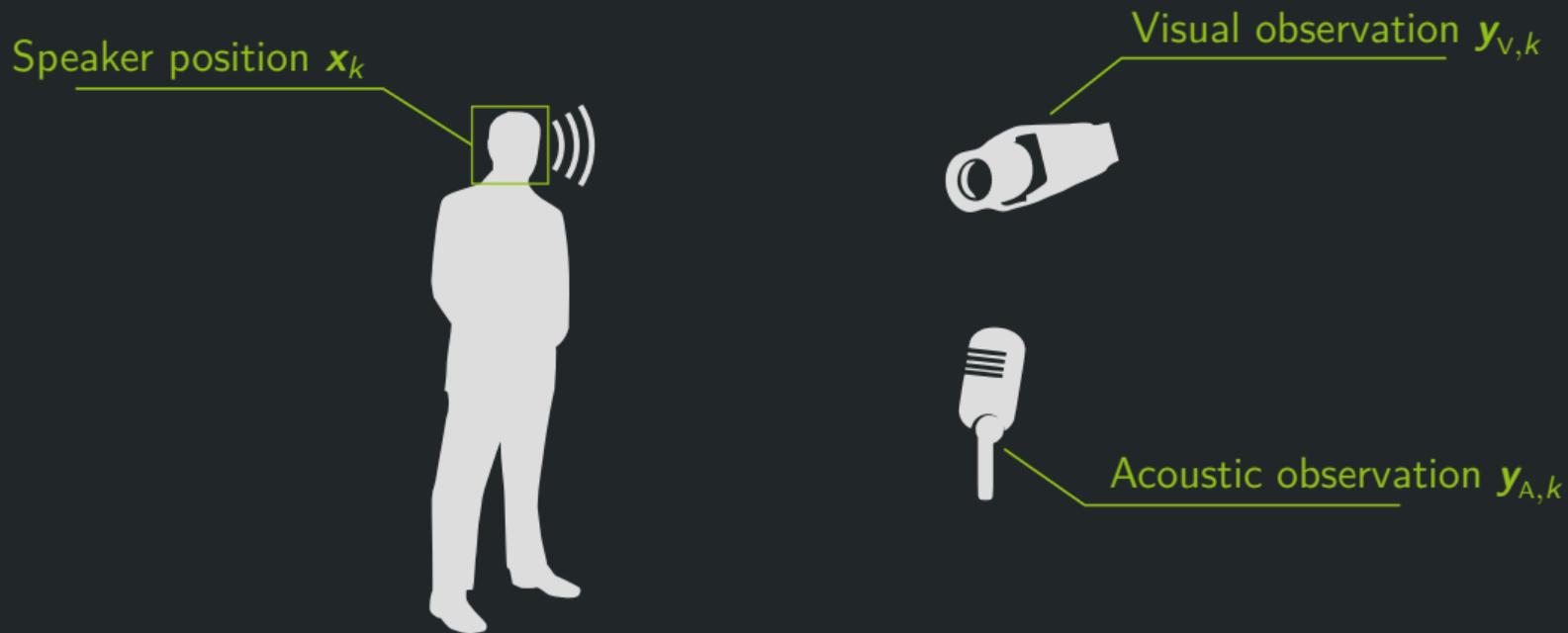


Audiovisual speaker tracking

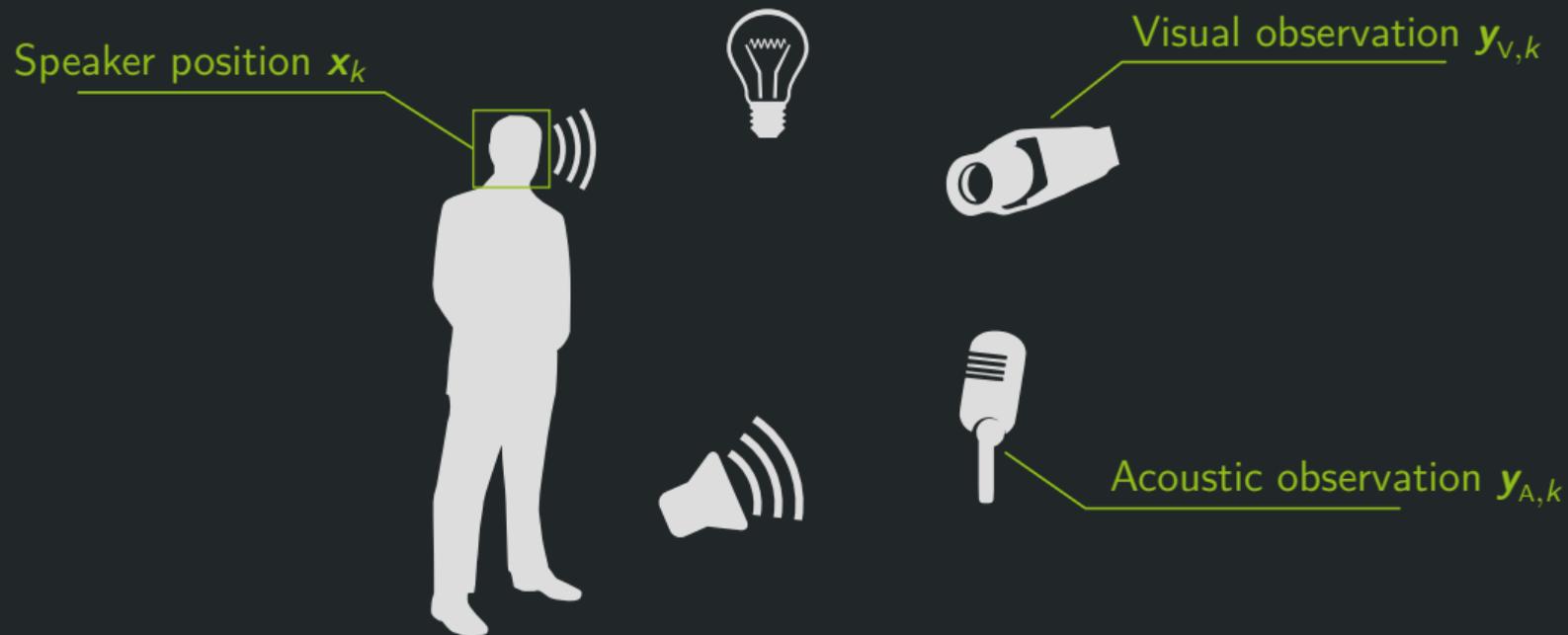
Speaker position \mathbf{x}_k



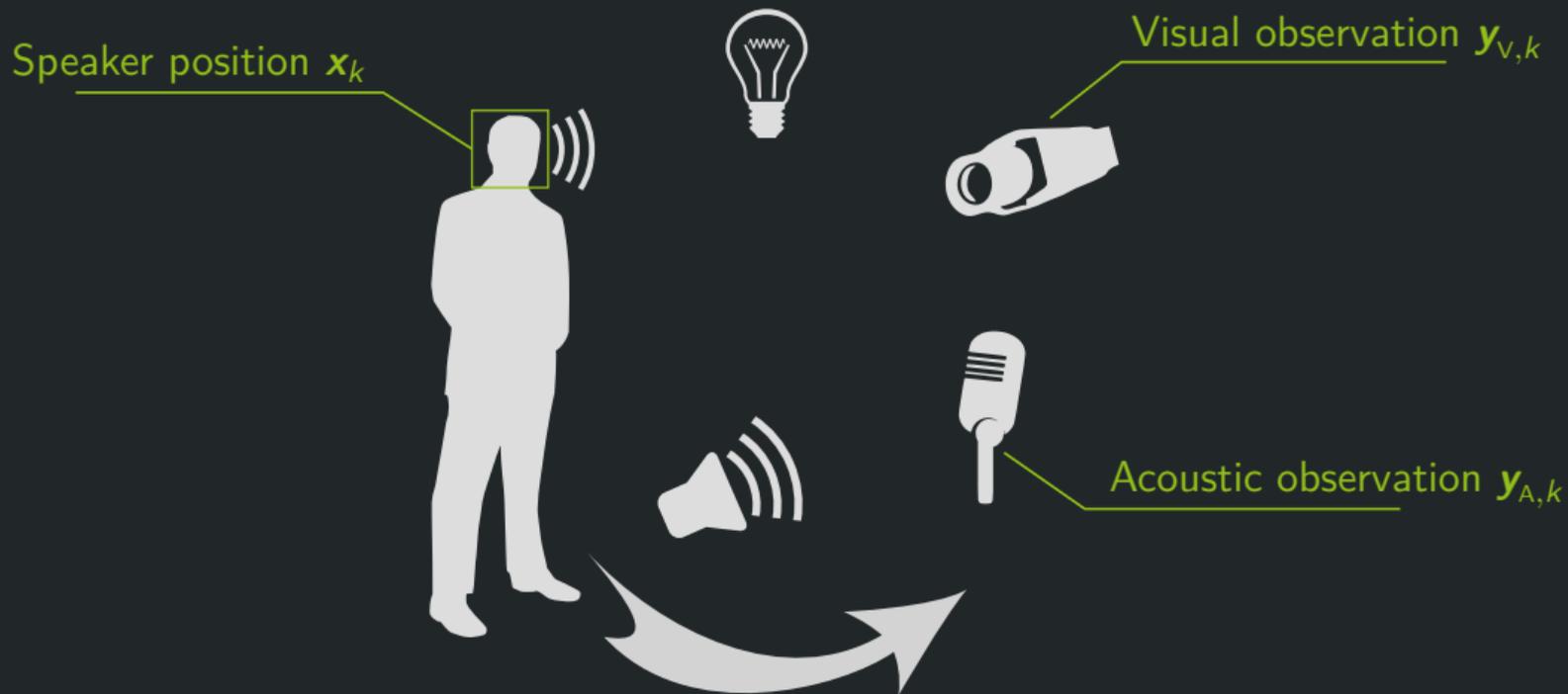
Audiovisual speaker tracking



Audiovisual speaker tracking



Audiovisual speaker tracking

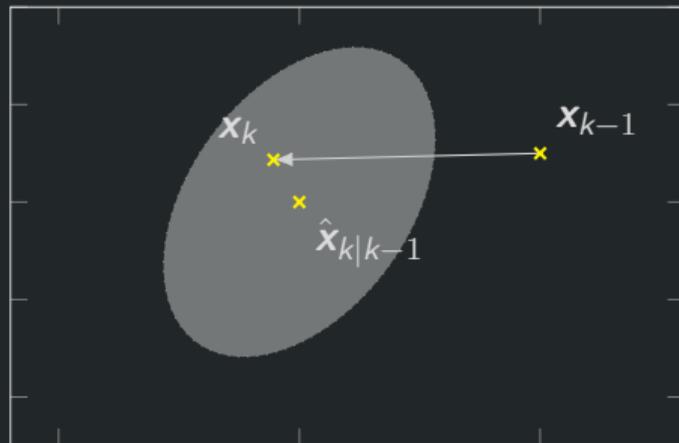


Audiovisual speaker tracking

Prediction step

System dynamics:

$$\mathbf{x}_k = f(\mathbf{x}_{k-1}) + \mathbf{v}_k, \quad \mathbf{v}_k = \mathcal{N}(\mathbf{0}, \mathbf{Q})$$

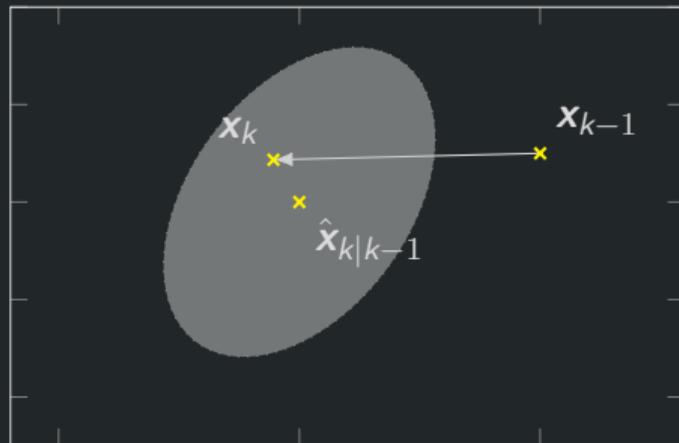


Audiovisual speaker tracking

Prediction step

System dynamics:

$$\mathbf{x}_k = f(\mathbf{x}_{k-1}) + \mathbf{v}_k, \quad \mathbf{v}_k = \mathcal{N}(\mathbf{0}, \mathbf{Q})$$



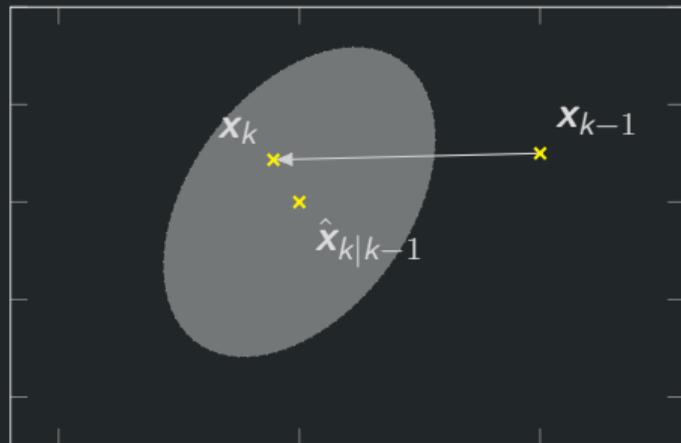
$$p(\mathbf{x}_k | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) = \int p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) d\mathbf{x}_{k-1}$$

Audiovisual speaker tracking

Prediction step

System dynamics:

$$\mathbf{x}_k = f(\mathbf{x}_{k-1}) + \mathbf{v}_k, \quad \mathbf{v}_k = \mathcal{N}(\mathbf{0}, \mathbf{Q})$$



$$p(\mathbf{x}_k | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) = \int \underbrace{p(\mathbf{x}_k | \mathbf{x}_{k-1})}_{\text{Dynamic model}} \underbrace{p(\mathbf{x}_{k-1} | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1})}_{\text{Prior}} d\mathbf{x}_{k-1}$$

Audiovisual speaker tracking

Observation

Observation model:

$$\mathbf{y}_k = \begin{bmatrix} \mathbf{y}_{A,k} & \mathbf{y}_{V,k} \end{bmatrix}^T = h(\mathbf{x}_k) + \mathbf{w}_k$$

$$\mathbf{w}_k = \mathcal{N}(\mathbf{0}, \mathbf{R}), \quad \mathbf{R} = \begin{bmatrix} \mathbf{R}_{AA} & \mathbf{R}_{AV} \\ \mathbf{R}_{VA} & \mathbf{R}_{VV} \end{bmatrix}$$



Audiovisual speaker tracking

Update step (standard Kalman filter)

Observation model:

$$\mathbf{y}_k = \begin{bmatrix} \mathbf{y}_{A,k} & \mathbf{y}_{V,k} \end{bmatrix}^T = h(\mathbf{x}_k) + \mathbf{w}_k$$

$$\mathbf{w}_k = \mathcal{N}(\mathbf{0}, \mathbf{R}), \quad \mathbf{R} = \begin{bmatrix} \mathbf{R}_{AA} & \mathbf{R}_{AV} \\ \mathbf{R}_{VA} & \mathbf{R}_{VV} \end{bmatrix}$$



Audiovisual speaker tracking

Update step (standard Kalman filter)

Observation model:

$$\mathbf{y}_k = \begin{bmatrix} \mathbf{y}_{A,k} & \mathbf{y}_{V,k} \end{bmatrix}^T = h(\mathbf{x}_k) + \mathbf{w}_k$$

$$\mathbf{w}_k = \mathcal{N}(\mathbf{0}, \mathbf{R}), \quad \mathbf{R} = \begin{bmatrix} \mathbf{R}_{AA} & \mathbf{R}_{AV} \\ \mathbf{R}_{VA} & \mathbf{R}_{VV} \end{bmatrix}$$



$$p(\mathbf{x}_k | \mathbf{Y}_{A,k}, \mathbf{Y}_{V,k}) \propto p(\mathbf{x}_k | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) p(\mathbf{y}_{A,k}, \mathbf{y}_{V,k} | \mathbf{x}_k)$$

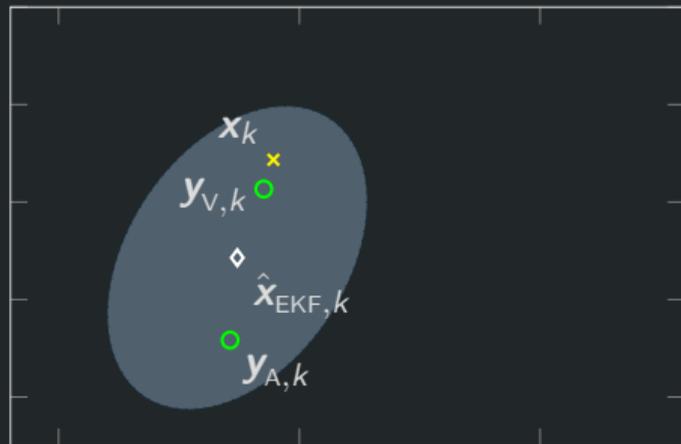
Audiovisual speaker tracking

Update step (standard Kalman filter)

Observation model:

$$\mathbf{y}_k = \begin{bmatrix} \mathbf{y}_{A,k} & \mathbf{y}_{V,k} \end{bmatrix}^T = h(\mathbf{x}_k) + \mathbf{w}_k$$

$$\mathbf{w}_k = \mathcal{N}(\mathbf{0}, \mathbf{R}), \quad \mathbf{R} = \begin{bmatrix} \mathbf{R}_{AA} & \mathbf{R}_{AV} \\ \mathbf{R}_{VA} & \mathbf{R}_{VV} \end{bmatrix}$$



$$p(\mathbf{x}_k | \mathbf{Y}_{A,k}, \mathbf{Y}_{V,k}) \propto p(\mathbf{x}_k | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) \underbrace{p(\mathbf{y}_{A,k}, \mathbf{y}_{V,k} | \mathbf{x}_k)}_{\text{Sensor model}}$$

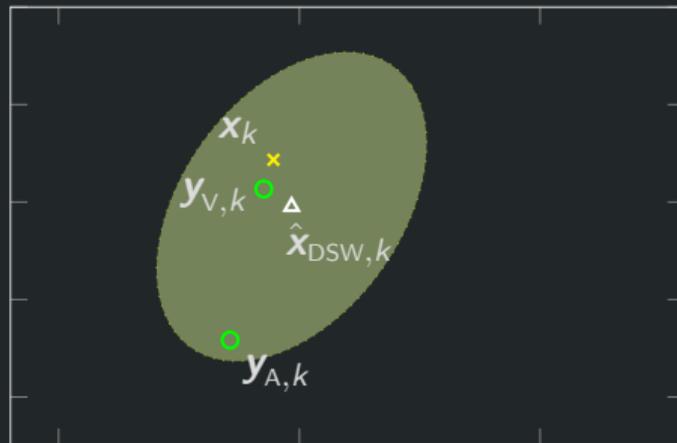
Audiovisual speaker tracking

Update step (Kalman filter with dynamic stream weights¹)

Observation model:

$$\mathbf{y}_{A,k} = h_A(\mathbf{x}_k) + \mathbf{w}_{A,k}, \quad \mathbf{w}_{A,k} = \mathcal{N}(\mathbf{0}, \mathbf{R}_{AA})$$

$$\mathbf{y}_{V,k} = h_V(\mathbf{x}_k) + \mathbf{w}_{V,k}, \quad \mathbf{w}_{V,k} = \mathcal{N}(\mathbf{0}, \mathbf{R}_{VV})$$



¹C. Schymura et al.: *Extending linear dynamical systems with dynamic stream weights for audiovisual speaker localization*, IWAENC, 2018

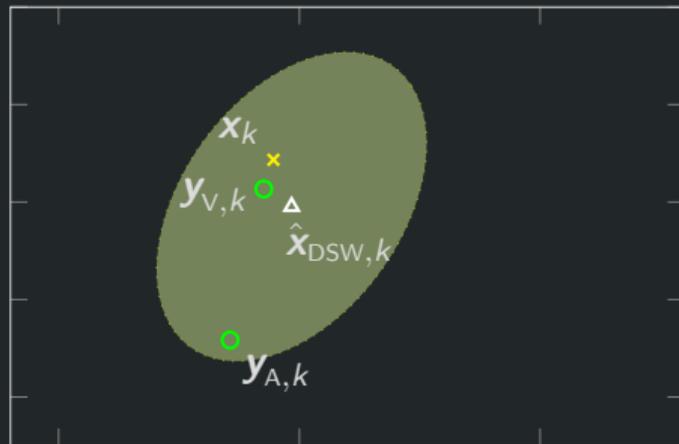
Audiovisual speaker tracking

Update step (Kalman filter with dynamic stream weights¹)

Observation model:

$$\mathbf{y}_{A,k} = h_A(\mathbf{x}_k) + \mathbf{w}_{A,k}, \quad \mathbf{w}_{A,k} = \mathcal{N}(\mathbf{0}, \mathbf{R}_{AA})$$

$$\mathbf{y}_{V,k} = h_V(\mathbf{x}_k) + \mathbf{w}_{V,k}, \quad \mathbf{w}_{V,k} = \mathcal{N}(\mathbf{0}, \mathbf{R}_{VV})$$



$$p(\mathbf{x}_k | \mathbf{Y}_{A,k}, \mathbf{Y}_{V,k}) \propto p(\mathbf{x}_k | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) \underbrace{p(\mathbf{y}_{A,k} | \mathbf{x}_k)^{\lambda_k}}_{\text{Acoustic model}} \underbrace{p(\mathbf{y}_{V,k} | \mathbf{x}_k)^{1-\lambda_k}}_{\text{Visual model}}$$

¹C. Schymura et al.: *Extending linear dynamical systems with dynamic stream weights for audiovisual speaker localization*, IWAENC, 2018

Inference

Extended Kalman filter approach: first-order Taylor series expansion

$$f(\mathbf{x}_{k-1}) \approx f(\hat{\mathbf{x}}_{k-1}) + \mathbf{F}(\hat{\mathbf{x}}_{k-1})(\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1})$$

Inference

Extended Kalman filter approach: first-order Taylor series expansion

$$f(\mathbf{x}_{k-1}) \approx f(\hat{\mathbf{x}}_{k-1}) + \mathbf{F}(\hat{\mathbf{x}}_{k-1})(\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1})$$

$$\Rightarrow p(\mathbf{x}_k | \mathbf{x}_{k-1}) = \mathcal{N}\left(\mathbf{x}_k | f(\hat{\mathbf{x}}_{k-1}) + \mathbf{F}(\hat{\mathbf{x}}_{k-1})(\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1}), \mathbf{Q}\right)$$

$$\Rightarrow p(\mathbf{x}_k | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) = \mathcal{N}\left(\mathbf{x}_k | \hat{\mathbf{x}}_{k-1}, \hat{\Sigma}_{k-1}\right)$$

Inference

Extended Kalman filter approach: first-order Taylor series expansion

$$f(\mathbf{x}_{k-1}) \approx f(\hat{\mathbf{x}}_{k-1}) + \mathbf{F}(\hat{\mathbf{x}}_{k-1})(\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1})$$

$$\Rightarrow p(\mathbf{x}_k | \mathbf{x}_{k-1}) = \mathcal{N}\left(\mathbf{x}_k | f(\hat{\mathbf{x}}_{k-1}) + \mathbf{F}(\hat{\mathbf{x}}_{k-1})(\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1}), \mathbf{Q}\right)$$

$$\Rightarrow p(\mathbf{x}_k | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) = \mathcal{N}\left(\mathbf{x}_k | \hat{\mathbf{x}}_{k-1}, \hat{\Sigma}_{k-1}\right)$$

Prediction step (identical to standard EKF)

$$\hat{\mathbf{x}}_{k|k-1} = f(\hat{\mathbf{x}}_{k-1})$$

$$\hat{\Sigma}_{k|k-1} = \mathbf{F}_{k-1} \hat{\Sigma}_{k-1} \mathbf{F}_{k-1}^T + \mathbf{Q}, \quad \mathbf{F}_{k-1} \equiv \mathbf{F}(\hat{\mathbf{x}}_{k-1}) = \left. \frac{\partial f(\mathbf{x}_{k-1})}{\partial \mathbf{x}_{k-1}} \right|_{\mathbf{x}_{k-1} = \hat{\mathbf{x}}_{k-1}}$$

Inference

Extended Kalman filter approach: first-order Taylor series expansion

$$h_{\{A,V\}}(\mathbf{x}_k) \approx h_{\{A,V\}}(\hat{\mathbf{x}}_k) + \mathbf{H}_{\{A,V\},k}(\mathbf{x}_k - \hat{\mathbf{x}}_k), \quad \mathbf{H}_{\{A,V\},k} \equiv \left. \frac{\partial h_{\{A,V\}}(\mathbf{x}_k)}{\partial \mathbf{x}_k} \right|_{\mathbf{x}_k = \hat{\mathbf{x}}_k}$$

Inference

Extended Kalman filter approach: first-order Taylor series expansion

$$h_{\{A,V\}}(\mathbf{x}_k) \approx h_{\{A,V\}}(\hat{\mathbf{x}}_k) + \mathbf{H}_{\{A,V\},k}(\mathbf{x}_k - \hat{\mathbf{x}}_k), \quad \mathbf{H}_{\{A,V\},k} \equiv \left. \frac{\partial h_{\{A,V\}}(\mathbf{x}_k)}{\partial \mathbf{x}_k} \right|_{\mathbf{x}_k = \hat{\mathbf{x}}_k}$$
$$\Rightarrow p(\mathbf{y}_{\{A,V\},k} | \mathbf{x}_k) = \mathcal{N}(\mathbf{y}_{\{A,V\},k} | h_{\{A,V\}}(\hat{\mathbf{x}}_k) + \mathbf{H}_{\{A,V\},k}(\mathbf{x}_k - \hat{\mathbf{x}}_k), \mathbf{R}_{\{A,V\}})$$

Inference

Extended Kalman filter approach: first-order Taylor series expansion

$$h_{\{A,V\}}(\mathbf{x}_k) \approx h_{\{A,V\}}(\hat{\mathbf{x}}_k) + \mathbf{H}_{\{A,V\},k}(\mathbf{x}_k - \hat{\mathbf{x}}_k), \quad \mathbf{H}_{\{A,V\},k} \equiv \left. \frac{\partial h_{\{A,V\}}(\mathbf{x}_k)}{\partial \mathbf{x}_k} \right|_{\mathbf{x}_k = \hat{\mathbf{x}}_k}$$
$$\Rightarrow p(\mathbf{y}_{\{A,V\},k} | \mathbf{x}_k) = \mathcal{N}\left(\mathbf{y}_{\{A,V\},k}, \mid h_{\{A,V\}}(\hat{\mathbf{x}}_k) + \mathbf{H}_{\{A,V\},k}(\mathbf{x}_k - \hat{\mathbf{x}}_k), \mathbf{R}_{\{A,V\}}\right)$$

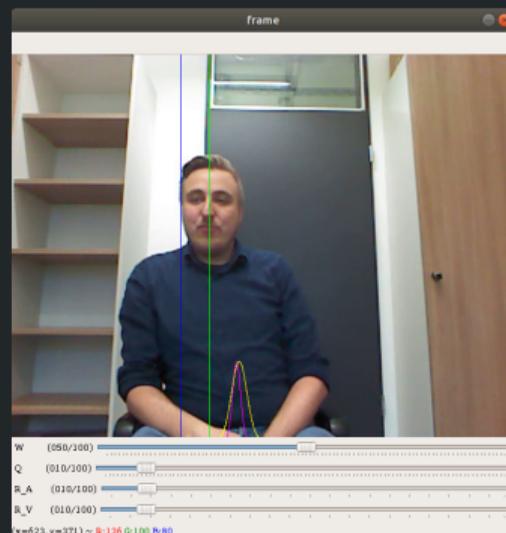
Update step

$$\begin{bmatrix} \mathbf{K}_{A,k}^T \\ \mathbf{K}_{V,k}^T \end{bmatrix} = \begin{bmatrix} \mathbf{R}_A + \lambda_k \mathbf{H}_{A,k} \hat{\Sigma}_{k|k-1} \mathbf{H}_{A,k}^T & (1 - \lambda_k) \mathbf{H}_{A,k} \hat{\Sigma}_{k|k-1} \mathbf{H}_{V,k}^T \\ \lambda_k \mathbf{H}_{V,k} \hat{\Sigma}_{k|k-1} \mathbf{H}_{A,k}^T & \mathbf{R}_V + (1 - \lambda_k) \mathbf{H}_{V,k} \hat{\Sigma}_{k|k-1} \mathbf{H}_{V,k}^T \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{H}_{A,k} \\ \mathbf{H}_{V,k} \end{bmatrix} \hat{\Sigma}_{k|k-1}$$
$$\hat{\mathbf{x}}_k = \hat{\mathbf{x}}_{k|k-1} + \lambda_k \mathbf{K}_{A,k} (\mathbf{y}_{A,k} - h_A(\hat{\mathbf{x}}_k)) + (1 - \lambda_k) \mathbf{K}_{V,k} (\mathbf{y}_{V,k} - h_V(\hat{\mathbf{x}}_k))$$
$$\hat{\Sigma}_{k|k-1} = \left(\mathbf{I} - \lambda_k \mathbf{K}_{A,k} \mathbf{H}_{A,k} - (1 - \lambda_k) \mathbf{K}_{V,k} \mathbf{H}_{V,k} \right) \hat{\Sigma}_{k|k-1}$$

Evaluation I

Experimental setup

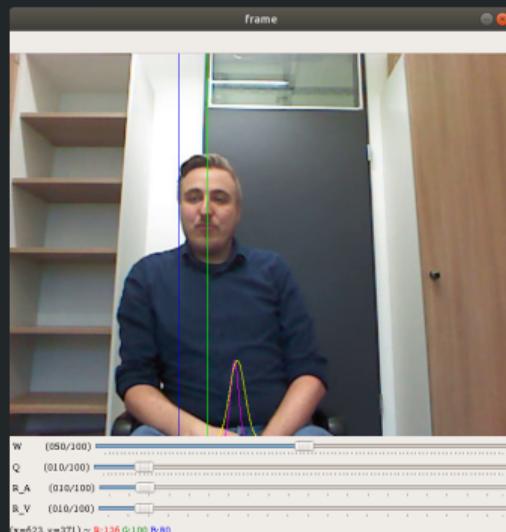
- ▶ KAVTraC audiovisual dataset, recorded in an office room at RUB using the Kinect sensor (7 speakers, $T_{60} \approx 350$ ms, 35 min. duration).



Evaluation I

Experimental setup

- ▶ KAVTraC audiovisual dataset, recorded in an office room at RUB using the Kinect sensor (7 speakers, $T_{60} \approx 350$ ms, 35 min. duration).
- ▶ Constant velocity linear dynamics model and nonlinear rotating vector observation models.
- ▶ DSW-EKF uses Dirichlet-prior oracle DSWs².

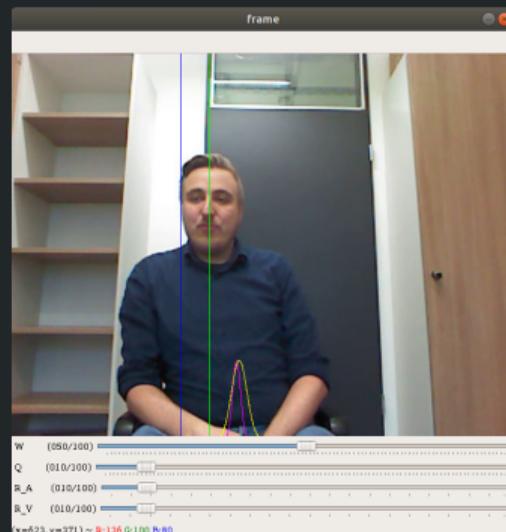


²C. Schymura et al.: *Audiovisual speaker tracking using nonlinear dynamical systems with dynamic stream weights*, arXiv, 2019

Evaluation I

Experimental setup

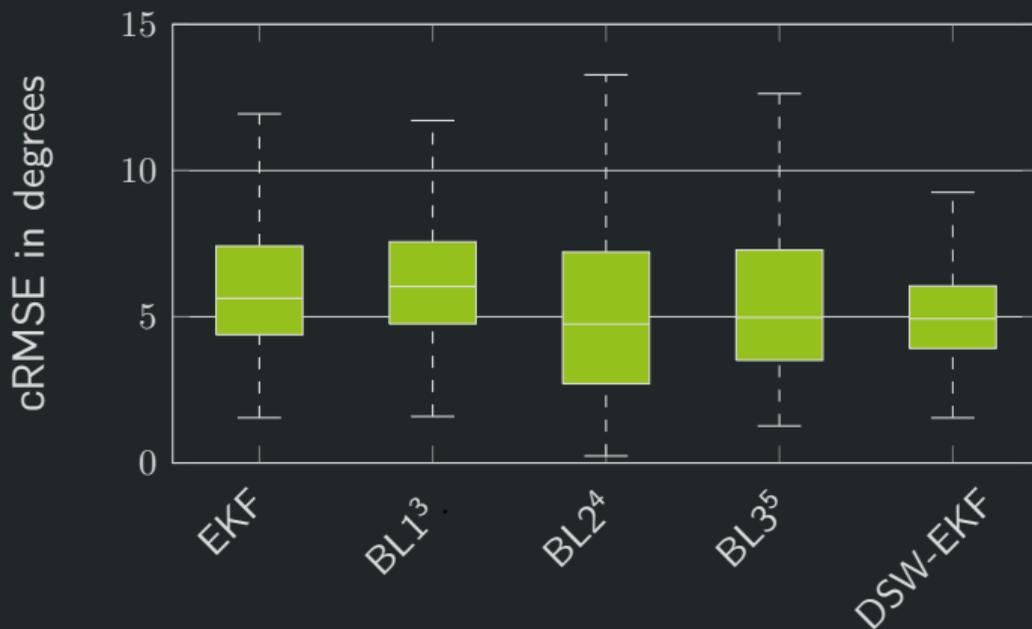
- ▶ KAVTraC audiovisual dataset, recorded in an office room at RUB using the Kinect sensor (7 speakers, $T_{60} \approx 350$ ms, 35 min. duration).
- ▶ Constant velocity linear dynamics model and nonlinear rotating vector observation models.
- ▶ DSW-EKF uses Dirichlet-prior oracle DSWs².
- ▶ Four baseline systems: standard EKF, one KF-based and two particle filter-based systems.
- ▶ Leave-one-out cross-validation paradigm.



²C. Schymura et al.: *Audiovisual speaker tracking using nonlinear dynamical systems with dynamic stream weights*, arXiv, 2019

Evaluation I

Results



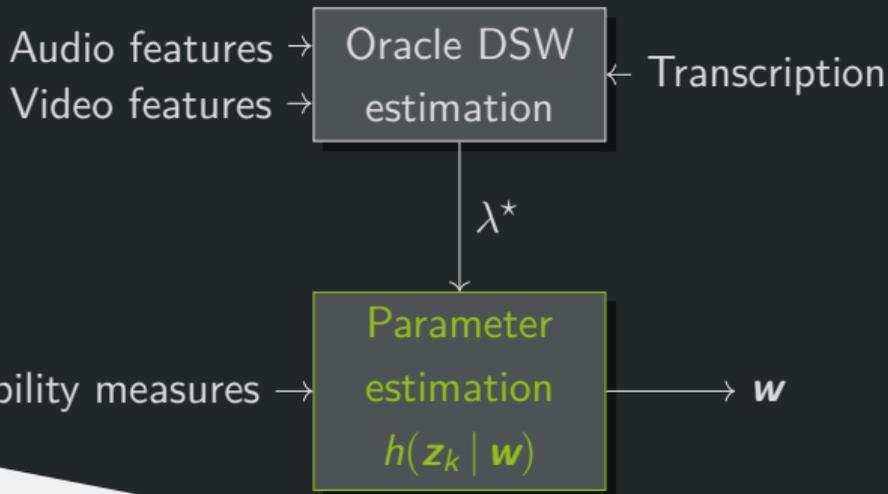
³T. Gehrig et al.: *Kalman filters for audio-video source localization*, WASPAA, 2005

⁴S. Gerlach et al.: *2D audio-visual localization in home environments using a particle filter*, ITG Symp., 2012

⁵X. Qian et al.: *3D audio-visual speaker tracking with an adaptive particle filter*, ICASSP, 2017

Learning dynamic stream weights

Standard approach: Supervised training with oracle dynamic stream weights



Learning Dynamic Stream Weights For Coupled-HMM-Based Audio-Visual Speech Recognition
Ahmed Hussien Abdelaziz, Student Member, IEEE, Steffen Zeiler, and Dorothea Kolossa, Senior Member, IEEE

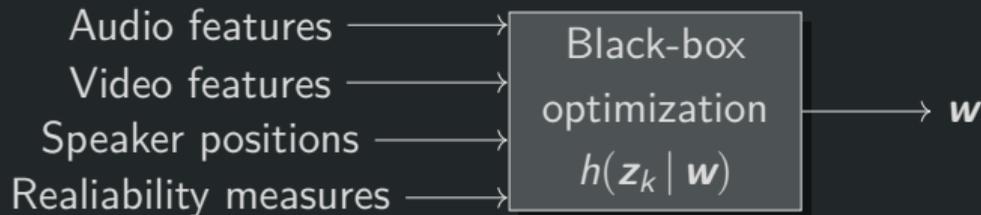
EXTENDING LINEAR DYNAMICAL SYSTEMS WITH DYNAMIC STREAM WEIGHTS FOR AUDIOVISUAL SPEAKER LOCALIZATION
Christopher Schymura, Tobias Isenberg and Dorothea Kolossa
Institute of Communication Acoustics, Ruhr-Universität Bochum, Germany

ABSTRACT
An important aspect of audiovisual speaker localization is the appropriate fusion of acoustic and visual observations based on their time-varying reliability. In this study, a framework which incorporates dynamic stream weights into the well-known linear dynamic stream weight estimator is proposed to consistently fuse audio-visual ASR system outputs.

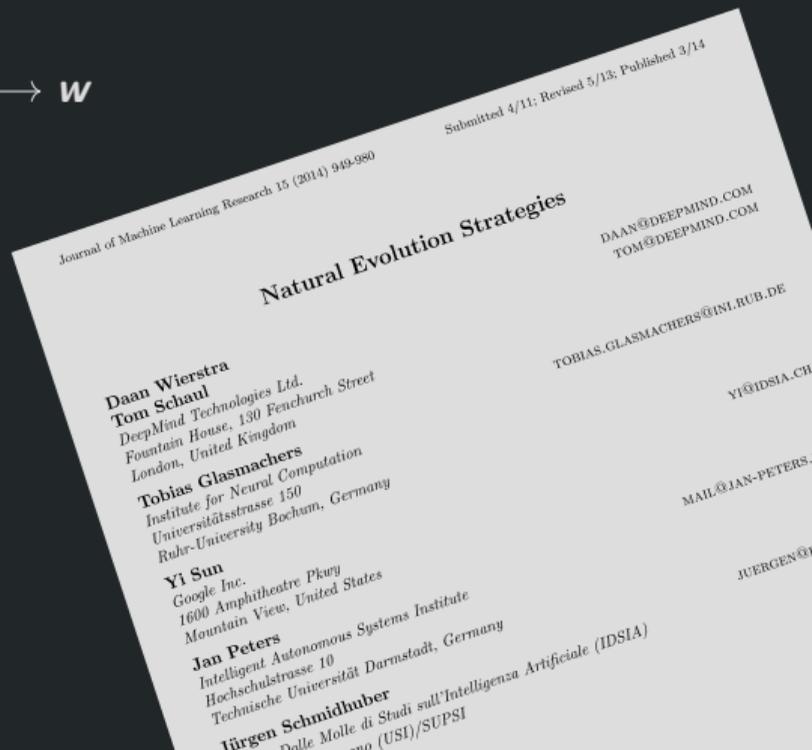
INDEX TERMS—Audio-visual speech recognition, dynamic stream weights, HMM-based reliability measure, speaker localization, supervised training.

Learning dynamic stream weights

Proposed approach: Training with natural evolution strategies



- ▶ No oracle information required.
- ▶ Flexible choice of loss/fitness function.



Learning dynamic stream weights

Training procedure⁶

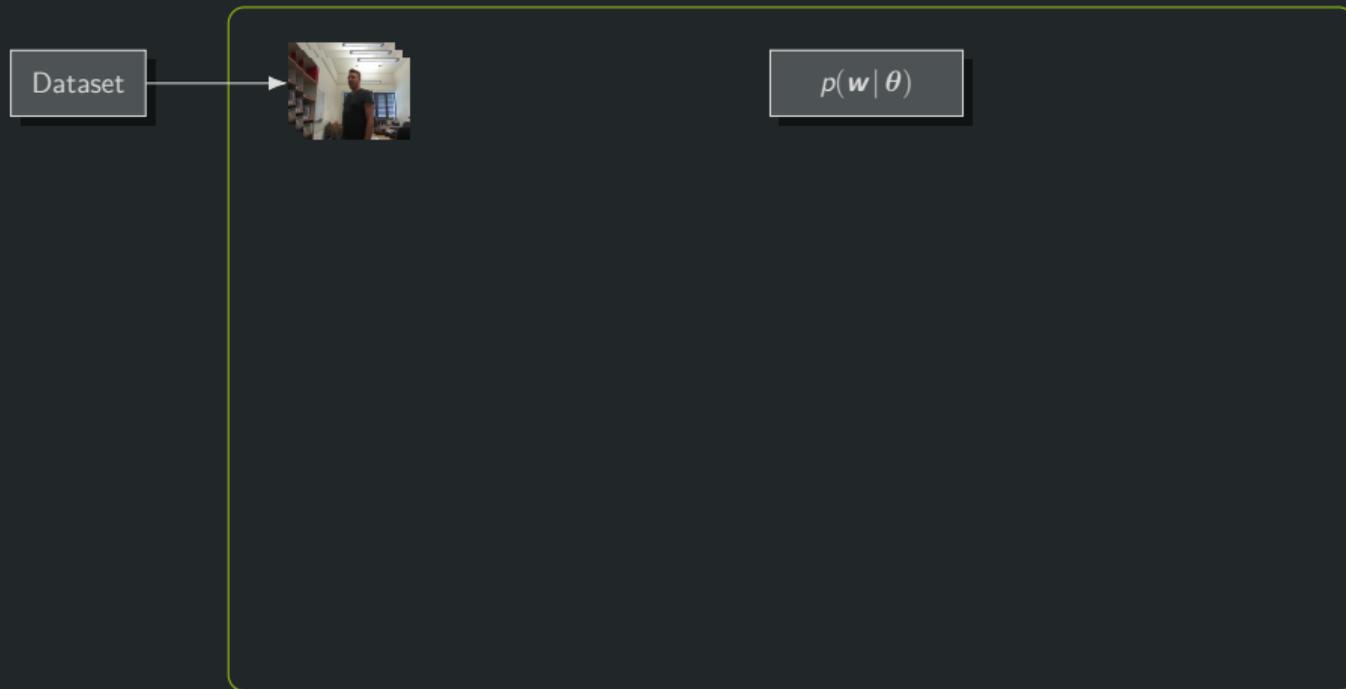
Dataset



⁶C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

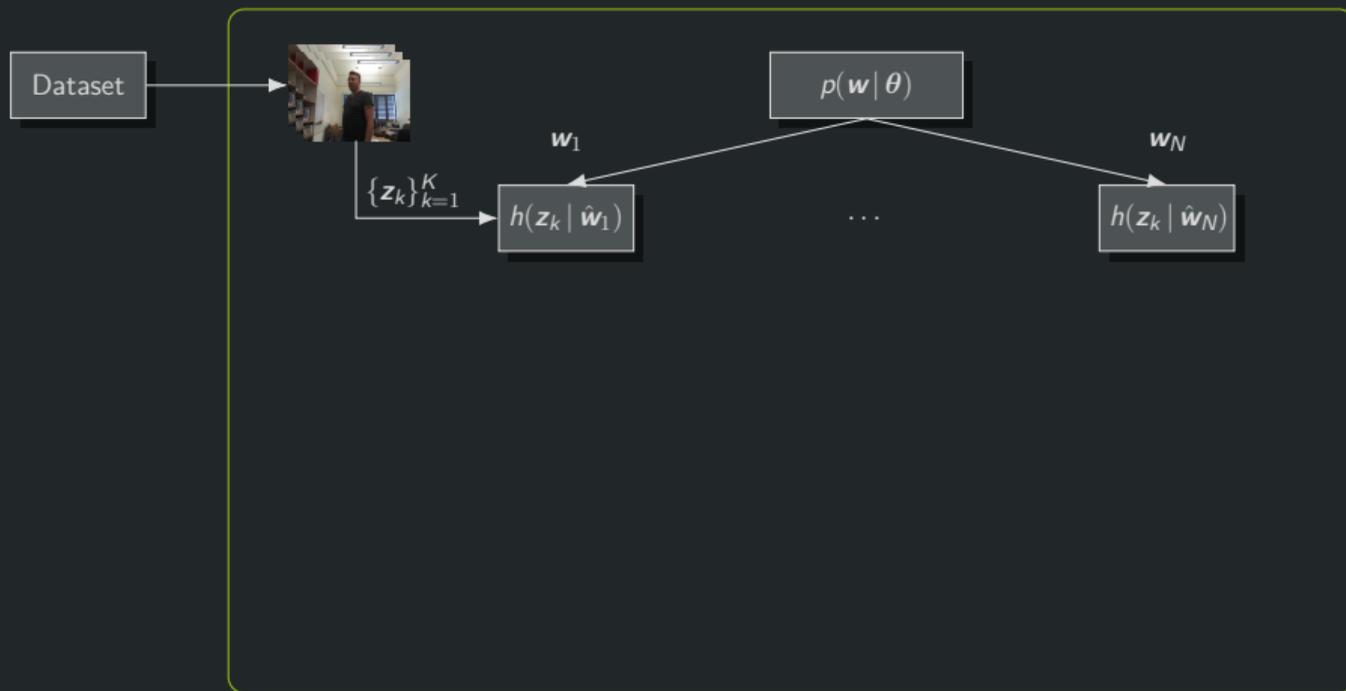
Training procedure⁶



⁶C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

Training procedure⁶



⁶C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

Training procedure⁶



⁶C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

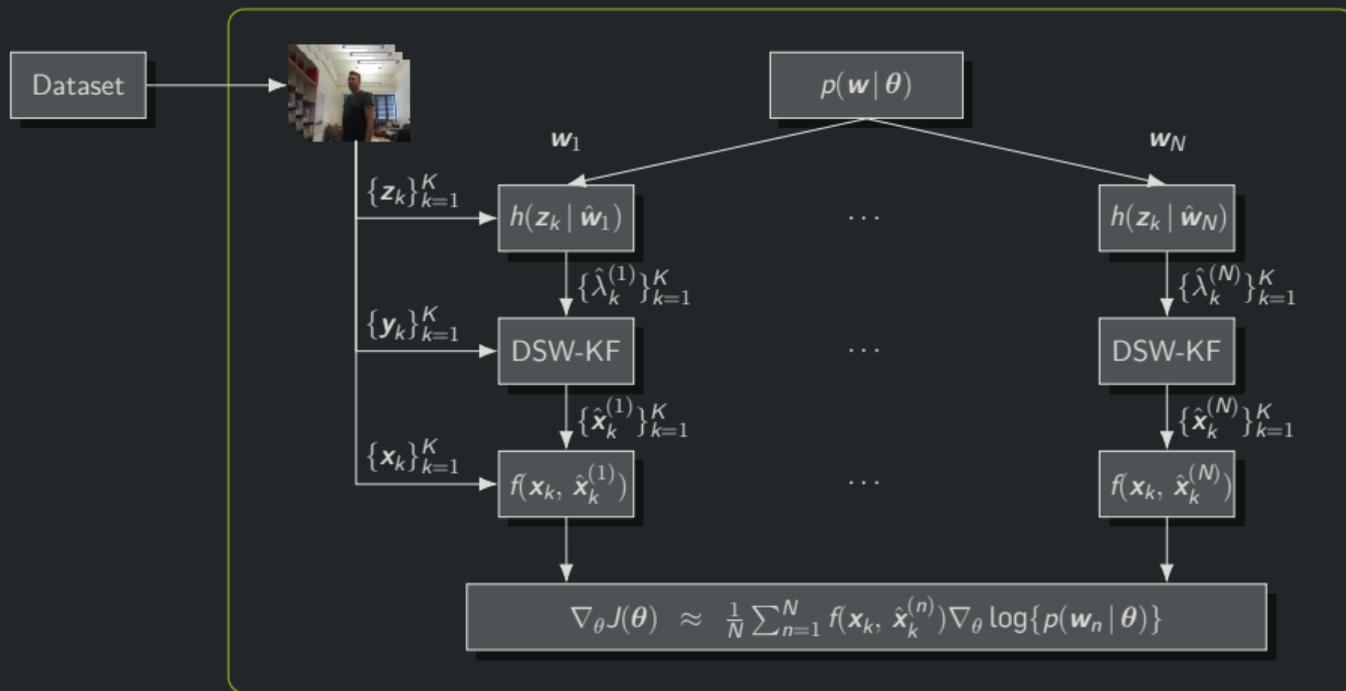
Training procedure⁶



⁶C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

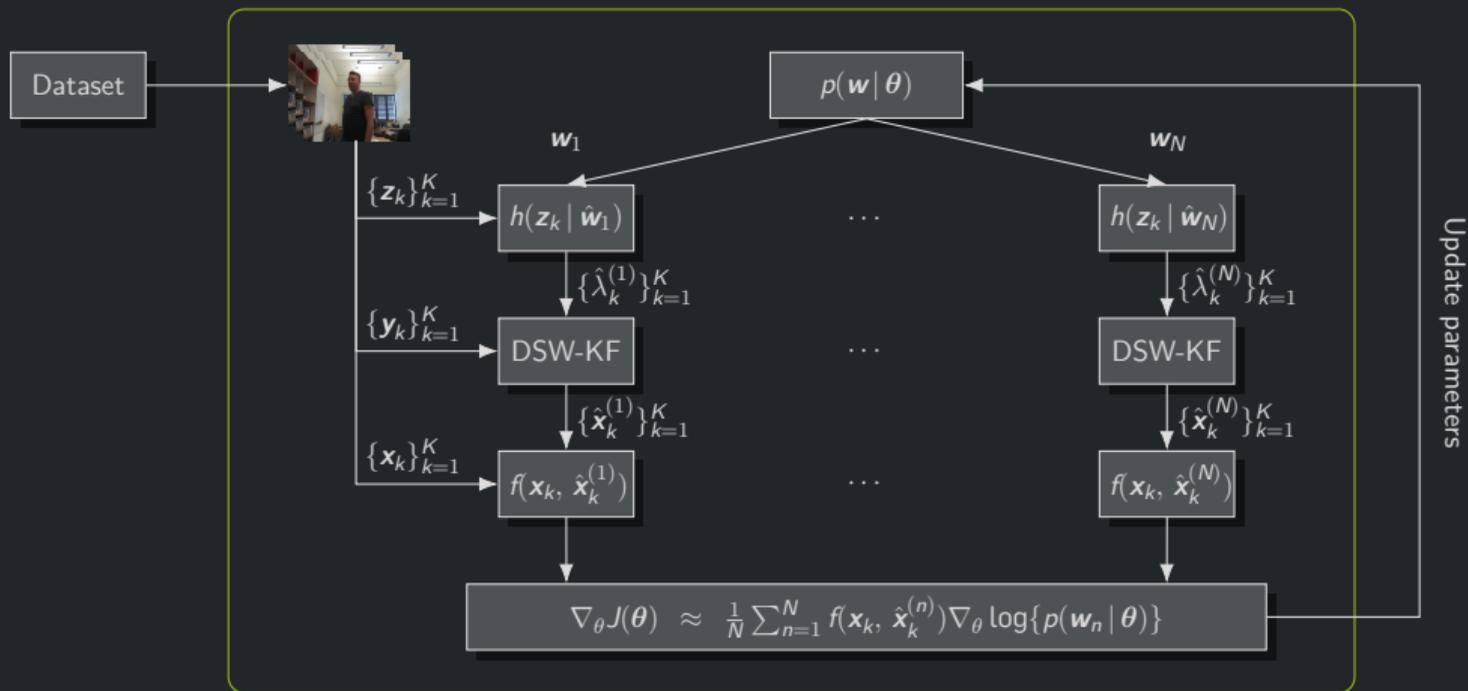
Training procedure⁶



⁶C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

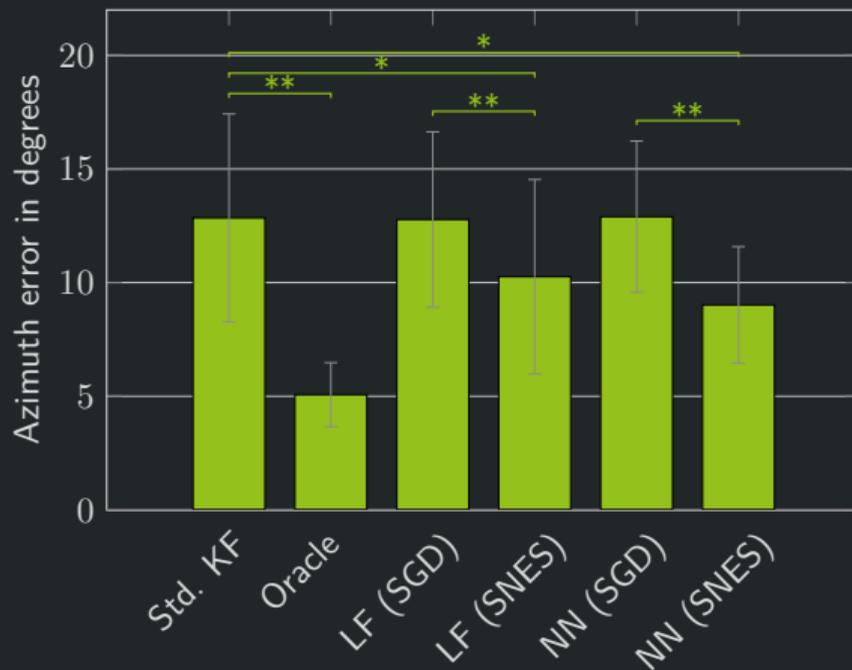
Training procedure⁶



⁶C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Evaluation II

Results



Conclusions and outlook

- ▶ DSW-based audiovisual speaker tracking frameworks can be extended to cope with nonlinear systems.

Conclusions and outlook

- ▶ DSW-based audiovisual speaker tracking frameworks can be extended to cope with nonlinear systems.
- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches.

Conclusions and outlook

- ▶ DSW-based audiovisual speaker tracking frameworks can be extended to cope with nonlinear systems.
- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches.
- ▶ Ideas for future work:

Conclusions and outlook

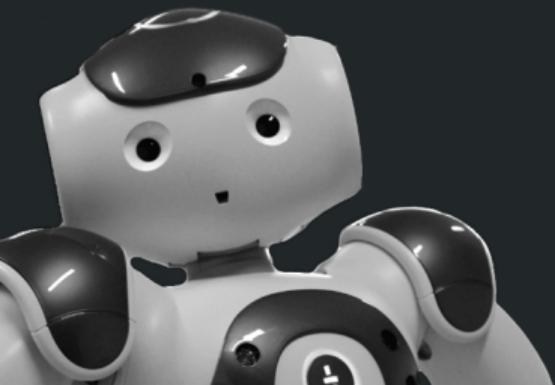
- ▶ DSW-based audiovisual speaker tracking frameworks can be extended to cope with nonlinear systems.
- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches.
- ▶ Ideas for future work:
 - ▶ Joint optimization of model and DSW estimation in a deep learning framework.

Conclusions and outlook

- ▶ DSW-based audiovisual speaker tracking frameworks can be extended to cope with nonlinear systems.
- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches.
- ▶ Ideas for future work:
 - ▶ Joint optimization of model and DSW estimation in a deep learning framework.
 - ▶ Extension to multi-speaker scenarios.

Conclusions and outlook

- ▶ DSW-based audiovisual speaker tracking frameworks can be extended to cope with nonlinear systems.
- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches.
- ▶ Ideas for future work:
 - ▶ Joint optimization of model and DSW estimation in a deep learning framework.
 - ▶ Extension to multi-speaker scenarios.



Thank you for your attention!