

# Learning Dynamic Stream Weights for Linear Dynamical Systems using Natural Evolution Strategies

ICASSP 2019

Christopher Schymura and Dorothea Kolossa

May 16th, 2019

RUHR  
UNIVERSITÄT  
BOCHUM

**RUB**

# Audiovisual speaker tracking



# Audiovisual speaker tracking

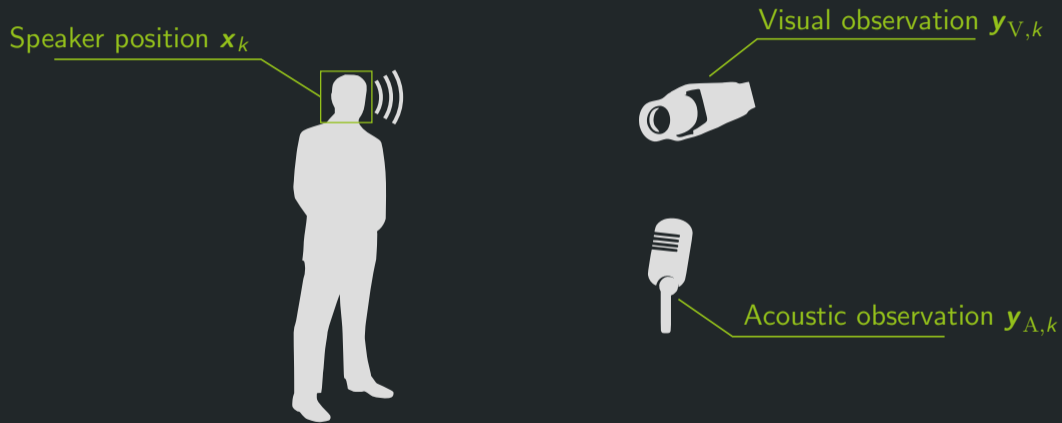


# Audiovisual speaker tracking

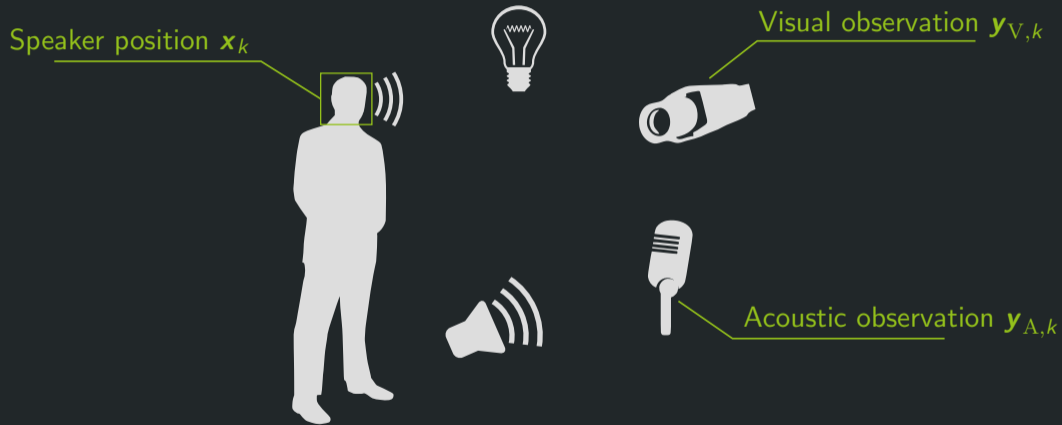
Speaker position  $\mathbf{x}_k$



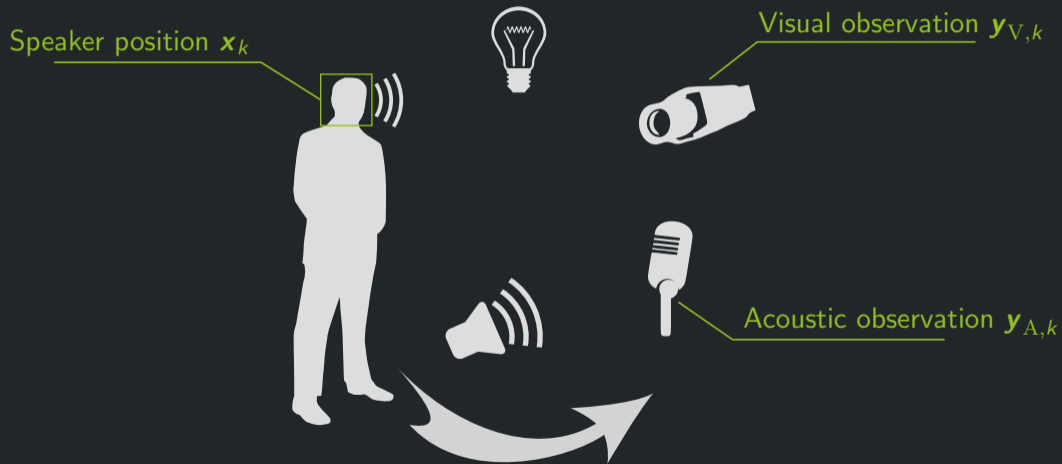
# Audiovisual speaker tracking



# Audiovisual speaker tracking



# Audiovisual speaker tracking

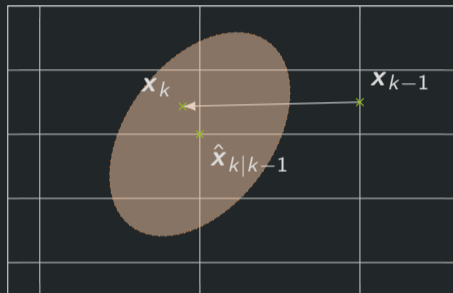


# Audiovisual speaker tracking

## Prediction step

System dynamics:

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



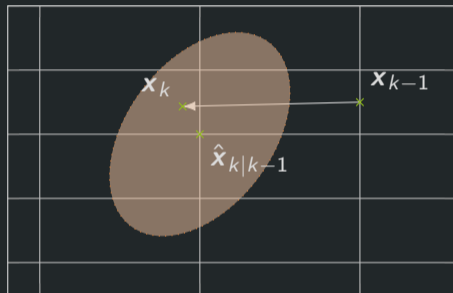


# Audiovisual speaker tracking

## Prediction step

System dynamics:

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



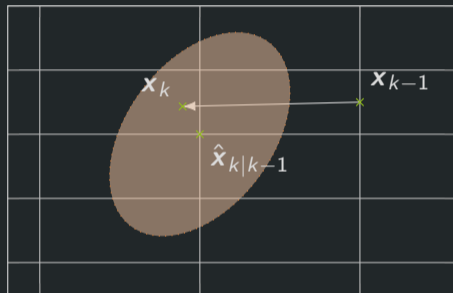
$$p(\mathbf{x}_k | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) = \int p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) d\mathbf{x}_{k-1}$$

# Audiovisual speaker tracking

## Prediction step

System dynamics:

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



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

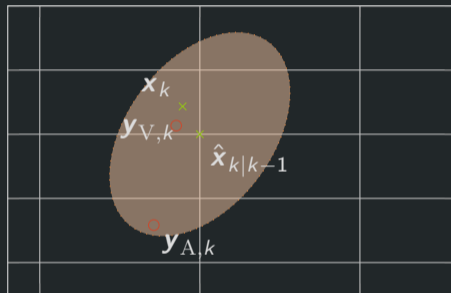
# Audiovisual speaker tracking

## Observation

Observation model:

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

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



# Audiovisual speaker tracking

## Update step (standard Kalman filter)

Observation model:

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

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



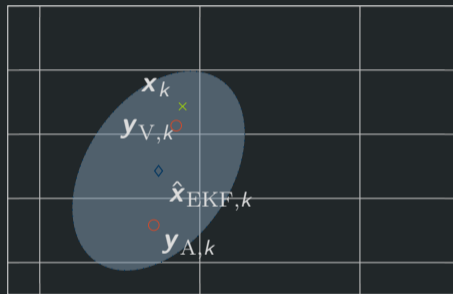
# Audiovisual speaker tracking

## Update step (standard Kalman filter)

Observation model:

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

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



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

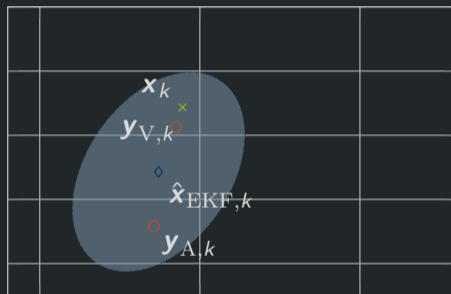
# Audiovisual speaker tracking

## Update step (standard Kalman filter)

Observation model:

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

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



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

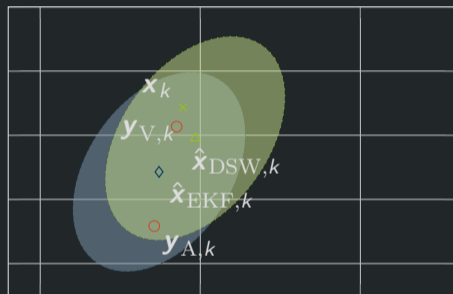
# Audiovisual speaker tracking

## Update step (Kalman filter with dynamic stream weights<sup>1</sup>)

Observation model:

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

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



<sup>1</sup>C. Schymura, T. Isenberg, D. Kolossa: *Extending Linear Dynamical Systems with Dynamic Stream Weights for Audiovisual Speaker*

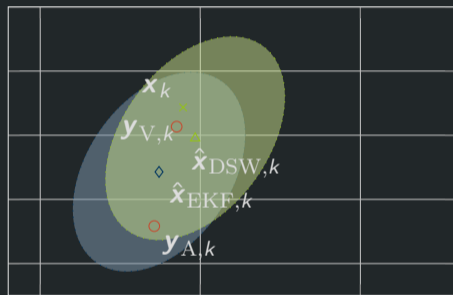
# Audiovisual speaker tracking

## Update step (Kalman filter with dynamic stream weights<sup>1</sup>)

Observation model:

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

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



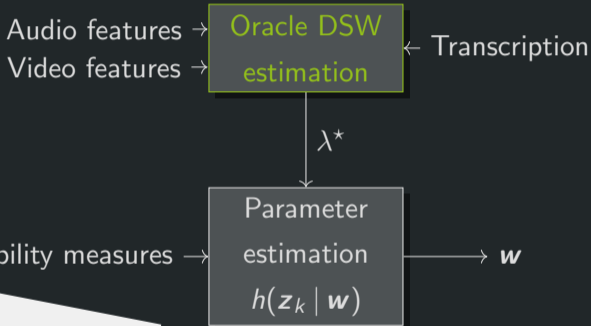
$$p(\mathbf{x}_k | \mathbf{Y}_{A,k}, \mathbf{Y}_{V,k}) \propto p(\mathbf{x}_k | \mathbf{Y}_{A,k-1}, \mathbf{Y}_{V,k-1}) \underbrace{p(\mathbf{y}_{A,k} | \mathbf{x}_k)^{\lambda_k}}_{\text{Acoustic model}} \underbrace{p(\mathbf{y}_{V,k} | \mathbf{x}_k)^{1-\lambda_k}}_{\text{Visual model}}$$

<sup>1</sup>C. Schymura, T. Isenberg, D. Kolossa: *Extending Linear Dynamical Systems with Dynamic Stream Weights for Audiovisual Speaker*



# 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

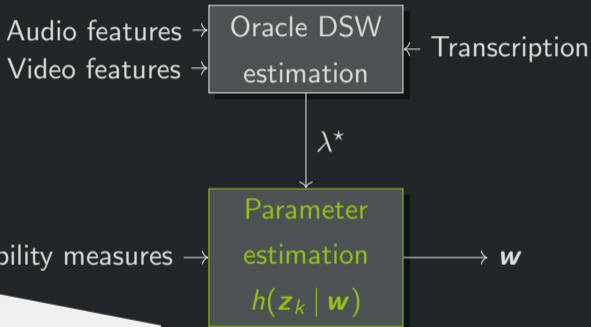
Abstract—With the increasing use of automatic speech recognition (ASR) systems in conjunction with the acoustic and visual information, dynamic stream weights for audio-visual systems enhance the recognition performance in noisy and dynamic environments according to their instantaneous reliability and systematically achieved. In this paper, we present a complete framework for learning dynamic stream weights for audio-visual systems based on coupled hidden Markov models (HMMs). We consider using dynamic stream weight estimators to map audio-visual stream weight estimator reliability measure features to audio-visual stream weight. We estimate these weights using an expectation-maximization (EM) algorithm. The 10-dimensional feature vector is used to estimate the reliability measure. During decoding, the reliability measure is used to blindly estimate the dynamic stream weights. The proposed framework significantly improves the performance of the audio-visual ASR system in noisy and dynamic environments.

**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 dynamical system (LDS) framework is proposed to consider the time-varying reliability of the visual observations. The visual observations relate to the position of a speaker in the scene. The visual sensor observations are modeled as a state-space model. The proposed framework is evaluated using the state-of-the-art stream weight estimator. The results show that the proposed framework significantly improves the performance of the audio-visual ASR system in noisy and dynamic environments.

# 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

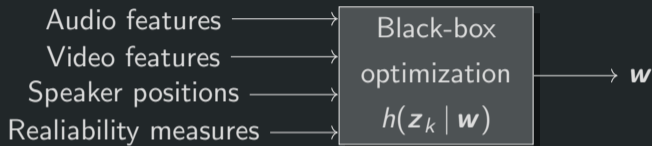
Abstract—With the increasing use of automatic speech recognition (ASR) systems in conjunction with the acoustic and visual information, dynamic stream weights for audio-visual systems enhance the recognition performance in noisy and dynamic environments. In this paper, we present a complete framework for learning dynamic stream weights for audio-visual systems. We consider using coupled hidden Markov models (HMMs) to model the joint audio-visual stream weight estimator. We estimate these weights using an expectation-maximization (EM) algorithm. We evaluate the proposed framework using the state-of-the-art audio-visual ASR system iSTRA.

**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

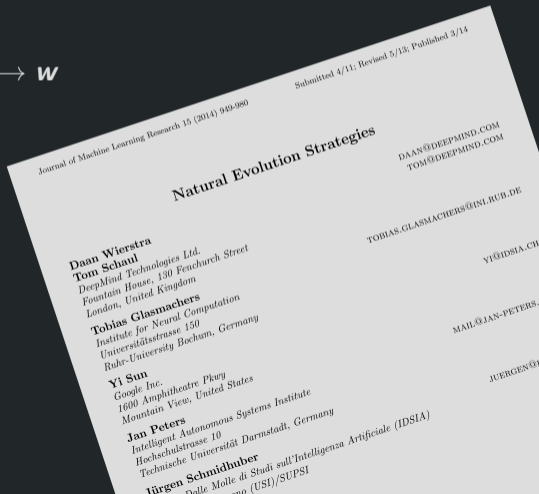
ABSTRACT  
An important aspect of audiovisual speaker localization is the appropriate fusion of acoustic and visual observations based on their time-varying reliability. In this study, a framework which incorporates dynamic stream weights into the well-known linear dynamical system (LDS) model is proposed to consider the reliability of the visual observations. The states represent the position of a speaker in the scene. The visual observations relate to measurements from multiple visual sensors.

# 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

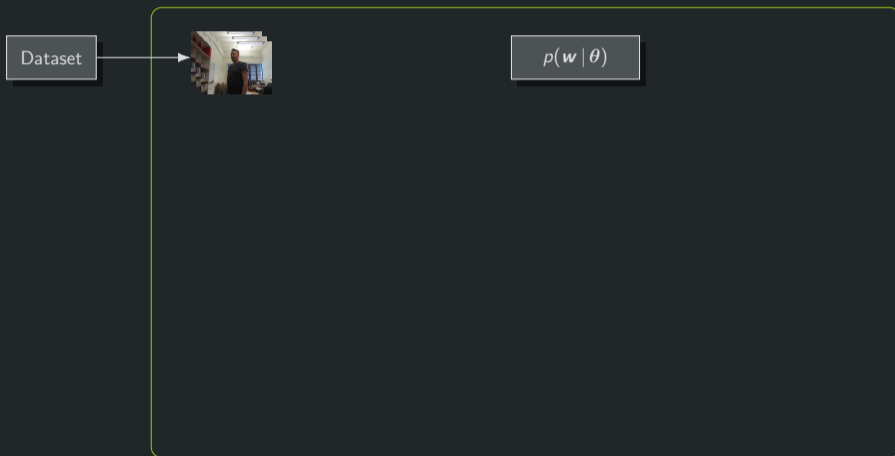
## Training procedure

Dataset



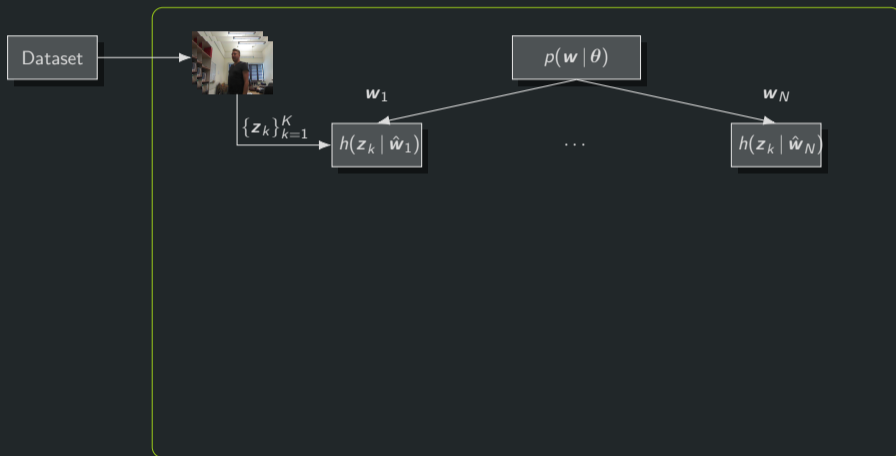
# Learning dynamic stream weights

## Training procedure



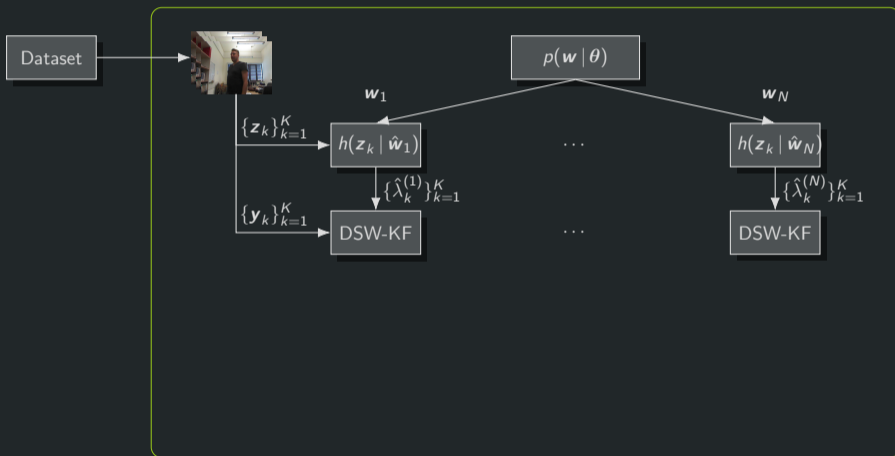
# Learning dynamic stream weights

## Training procedure



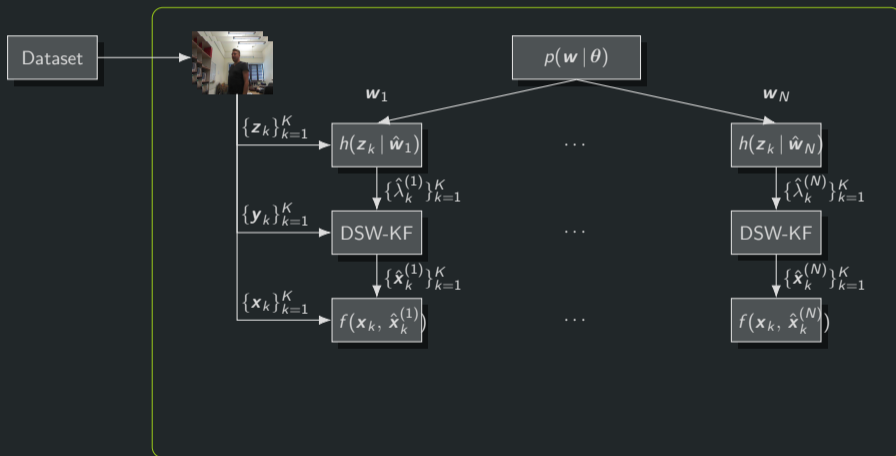
# Learning dynamic stream weights

## Training procedure



# Learning dynamic stream weights

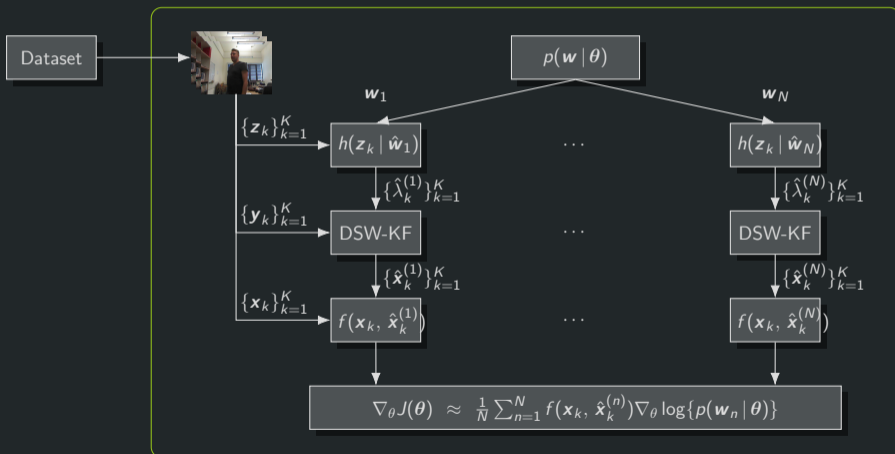
## Training procedure





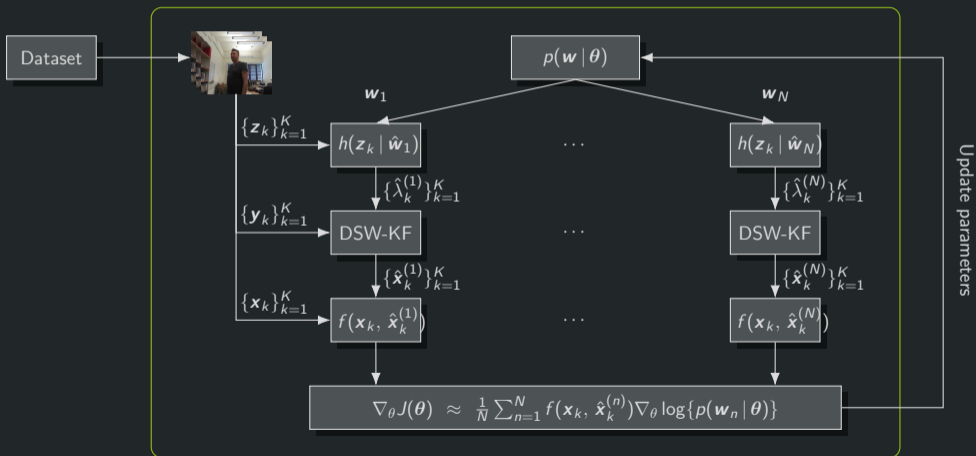
# Learning dynamic stream weights

## Training procedure



# Learning dynamic stream weights

## Training procedure



# Learning dynamic stream weights

## Implementation

- ▶ Reliability measures: instantaneous estimated a-priori SNR, acoustic and visual observation log-likelihoods<sup>2</sup>.

---

<sup>2</sup>A. H. Abdelaziz, S. Zeiler, D. Kolossa: *Learning Dynamic Stream Weights for Coupled-HMM-Based Audio-Visual Speech Recognition*, 2015

# Learning dynamic stream weights

## Implementation

- ▶ Reliability measures: instantaneous estimated a-priori SNR, acoustic and visual observation log-likelihoods<sup>2</sup>.
- ▶ Evaluation of two different DSW prediction models: logistic function and fully-connected feed-forward neural network.

---

<sup>2</sup>A. H. Abdelaziz, S. Zeiler, D. Kolossa: *Learning Dynamic Stream Weights for Coupled-HMM-Based Audio-Visual Speech Recognition*, 2015

# Learning dynamic stream weights

## Implementation

- ▶ Reliability measures: instantaneous estimated a-priori SNR, acoustic and visual observation log-likelihoods<sup>2</sup>.
- ▶ Evaluation of two different DSW prediction models: logistic function and fully-connected feed-forward neural network.
- ▶ Separable natural evolution strategies (sNES) as optimizer:

$$p(\mathbf{w} | \theta) = \mathcal{N}\left(\mathbf{w} | \mu_{\mathbf{w}}, \text{diag}(\sigma_{\mathbf{w}})\right)$$

---

<sup>2</sup>A. H. Abdelaziz, S. Zeiler, D. Kolossa: *Learning Dynamic Stream Weights for Coupled-HMM-Based Audio-Visual Speech Recognition*, 2015

# Learning dynamic stream weights

## Implementation

- ▶ Reliability measures: instantaneous estimated a-priori SNR, acoustic and visual observation log-likelihoods<sup>2</sup>.
- ▶ Evaluation of two different DSW prediction models: logistic function and fully-connected feed-forward neural network.

- ▶ Separable natural evolution strategies (sNES) as optimizer:

$$p(\mathbf{w} | \boldsymbol{\theta}) = \mathcal{N}\left(\mathbf{w} | \boldsymbol{\mu}_{\mathbf{w}}, \text{diag}(\boldsymbol{\sigma}_{\mathbf{w}})\right)$$

- ▶ Fitness function allowing direct optimization of instantaneous localization error:

$$f(\mathbf{w}) = -\frac{1}{M} \sum_{m=1}^M \frac{1}{K_m} \sum_{k=1}^{K_m} \left( \phi_k^{(m)} - \hat{\phi}_k^{(m)}(\mathbf{w}) \right)^2$$

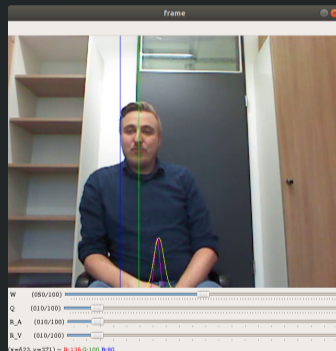
---

<sup>2</sup>A. H. Abdelaziz, S. Zeiler, D. Kolossa: *Learning Dynamic Stream Weights for Coupled-HMM-Based Audio-Visual Speech Recognition*, 2015

# Evaluation

## Experimental setup

- ▶ Front-end: DPD-MUSIC<sup>3</sup> for acoustic localization, Viola-Jones<sup>4</sup> algorithm for visual localization.



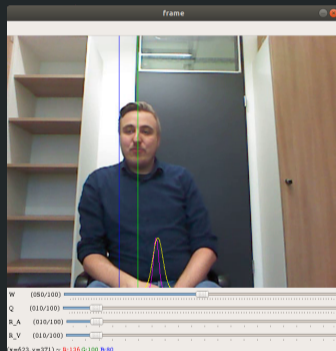
<sup>3</sup>Nadiri et al.: *Localization of multiple speakers under high reverberation using a spherical microphone array and the direct-path dominance test*, 2014

<sup>4</sup>P. Viola, M. Jones: *Rapid object detection using a boosted cascade of simple features*, 2001

# Evaluation

## Experimental setup

- ▶ Front-end: DPD-MUSIC<sup>3</sup> for acoustic localization, Viola-Jones<sup>4</sup> algorithm for visual localization.
- ▶ Dataset of audiovisual recordings in an office environment ( $T_{60} \approx 350$  ms) using the Kinect.



<sup>3</sup>Nadiri et al.: *Localization of multiple speakers under high reverberation using a spherical microphone array and the direct-path dominance test*, 2014

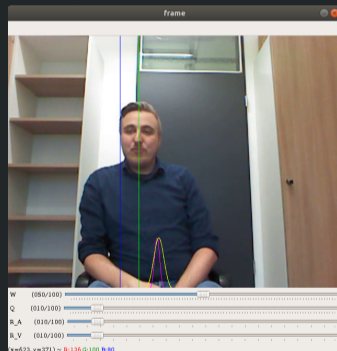
<sup>4</sup>P. Viola, M. Jones: *Rapid object detection using a boosted cascade of simple features*, 2001



# Evaluation

## Experimental setup

- ▶ Front-end: DPD-MUSIC<sup>3</sup> for acoustic localization, Viola-Jones<sup>4</sup> algorithm for visual localization.
- ▶ Dataset of audiovisual recordings in an office environment ( $T_{60} \approx 350$  ms) using the Kinect.
- ▶ Constant velocity dynamics model.



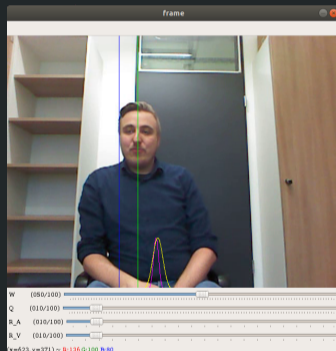
<sup>3</sup>Nadiri et al.: *Localization of multiple speakers under high reverberation using a spherical microphone array and the direct-path dominance test*, 2014

<sup>4</sup>P. Viola, M. Jones: *Rapid object detection using a boosted cascade of simple features*, 2001

# Evaluation

## Experimental setup

- ▶ Front-end: DPD-MUSIC<sup>3</sup> for acoustic localization, Viola-Jones<sup>4</sup> algorithm for visual localization.
- ▶ Dataset of audiovisual recordings in an office environment ( $T_{60} \approx 350$  ms) using the Kinect.
- ▶ Constant velocity dynamics model.
- ▶ Baseline: Stream weight prediction models trained on oracle DSWs with SGD (same architecture)

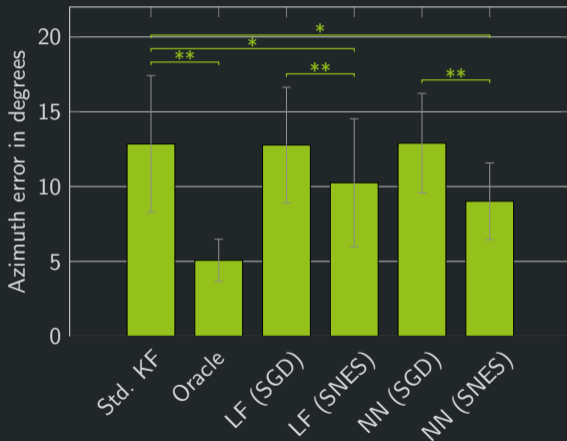


<sup>3</sup>Nadiri et al.: *Localization of multiple speakers under high reverberation using a spherical microphone array and the direct-path dominance test*, 2014

<sup>4</sup>P. Viola, M. Jones: *Rapid object detection using a boosted cascade of simple features*, 2001

# Evaluation

## Results



Statistical significance: \* for  $p < 0.05$  and \*\* for  $p < 0.01$

## Conclusions and outlook

- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches like NES (no oracle DSWs required).

## Conclusions and outlook

- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches like NES (no oracle DSWs required).
- ▶ Ideas for future work:

## Conclusions and outlook

- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches like NES (no oracle DSWs required).
- ▶ Ideas for future work:
  - ▶ Making the system trainable end-to-end.

## Conclusions and outlook

- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches like NES (no oracle DSWs required).
- ▶ Ideas for future work:
  - ▶ Making the system trainable end-to-end.
  - ▶ Joint optimization of DSW estimators and model parameters.

## Conclusions and outlook

- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches like NES (no oracle DSWs required).
- ▶ Ideas for future work:
  - ▶ Making the system trainable end-to-end.
  - ▶ Joint optimization of DSW estimators and model parameters.
  - ▶ Extension to multi-speaker scenarios.



## Conclusions and outlook

- ▶ A DSW-based audiovisual speaker tracking system can benefit from black-box optimization approaches like NES (no oracle DSWs required).
- ▶ Ideas for future work:
  - ▶ Making the system trainable end-to-end.
  - ▶ Joint optimization of DSW estimators and model parameters.
  - ▶ Extension to multi-speaker scenarios.