# 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

**RU**B

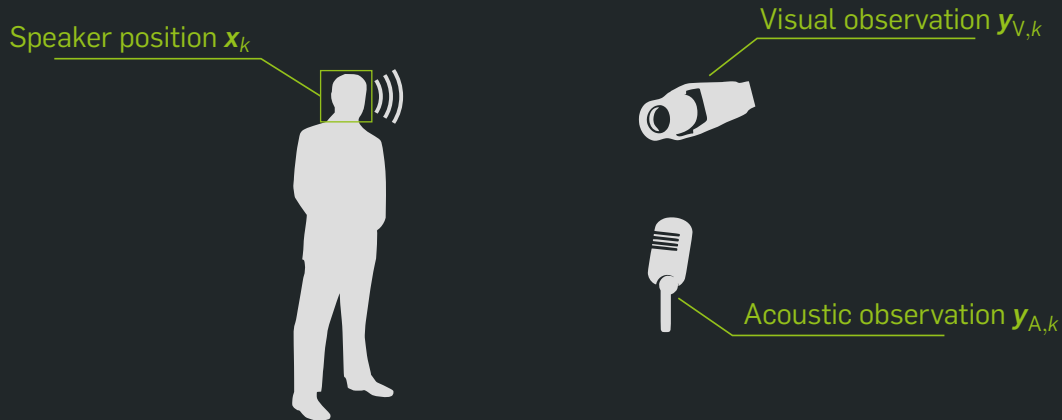# Audiovisual speaker tracking

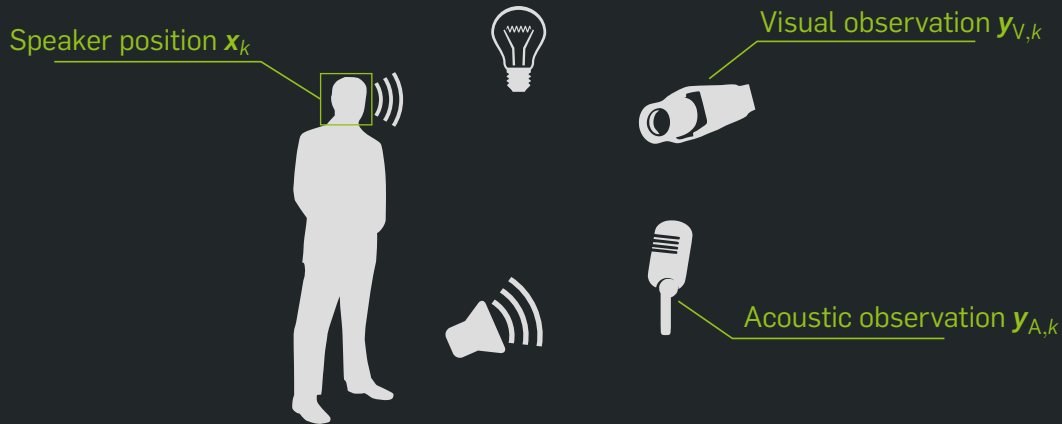# Audiovisual speaker tracking

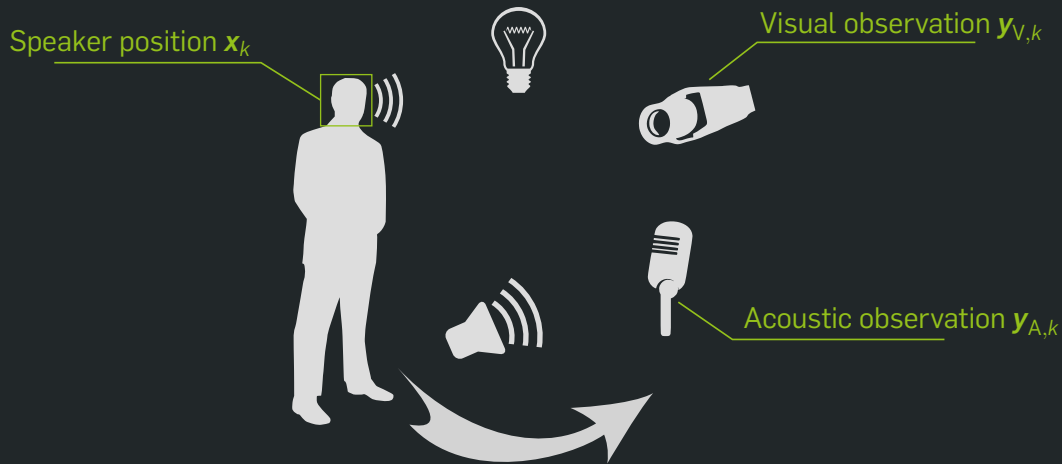# Audiovisual speaker tracking



Speaker position $\boldsymbol{x}_k$

# Audiovisual speaker tracking

# Audiovisual speaker tracking



Speaker position $\boldsymbol{x}_k$

Visual observation $\boldsymbol{y}_{V,k}$

Acoustic observation $\boldsymbol{y}_{A,k}$

# Audiovisual speaker tracking



Speaker position $\boldsymbol{x}_k$

Visual observation $\boldsymbol{y}_{\mathrm{V},k}$

Acoustic observation $\boldsymbol{y}_{\mathrm{A},k}$

# Audiovisual speaker tracking

**Prediction step**

System dynamics:

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

# Audiovisual speaker tracking

**Prediction step**

System dynamics:

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



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

# Audiovisual speaker tracking

**Prediction step**

System dynamics:

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



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

# Audiovisual speaker tracking

**Observation**

Observation model:

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

$$\boldsymbol{w}_k = \mathcal{N}(\boldsymbol{0}, \boldsymbol{R}), \quad \boldsymbol{R} = \begin{bmatrix} \boldsymbol{R}_{\mathsf{AA}} & \boldsymbol{R}_{\mathsf{AV}} \\ \boldsymbol{R}_{\mathsf{VA}} & \boldsymbol{R}_{\mathsf{VV}} \end{bmatrix}$$
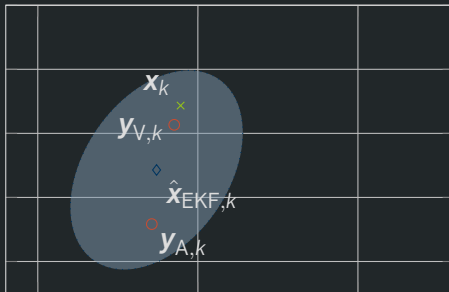
# Audiovisual speaker tracking

**Update step (standard Kalman filter)**

Observation model:

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

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

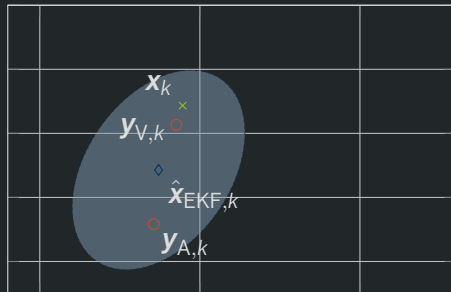# Audiovisual speaker tracking

**Update step (standard Kalman filter)**

Observation model:

$$y_k = \begin{bmatrix} y_{A,k} & y_{V,k} \end{bmatrix}^\top = Cx_k + w_k$$

$$w_k = \mathcal{N}(0, R), \quad R = \begin{bmatrix} R_{AA} & R_{AV} \\ R_{VA} & R_{VV} \end{bmatrix}$$



$$p(x_k \,|\, Y_{A,k}, Y_{V,k}) \propto p(x_k \,|\, Y_{A,k-1}, Y_{V,k-1})\, p(y_{A,k}, y_{V,k} \,|\, x_k)$$
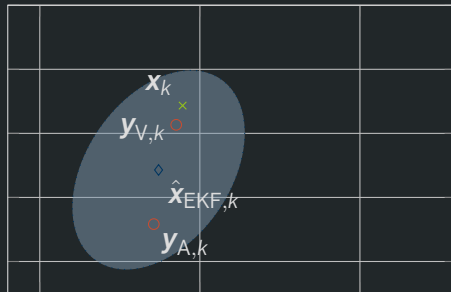
# Audiovisual speaker tracking

**Update step (standard Kalman filter)**

Observation model:

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

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



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

# Audiovisual speaker tracking

**Update step (Kalman filter with dynamic stream weights[1])**

Observation model:

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

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

[1] C. Schymura, T. Isenberg, D. Kolossa: *Extending Linear Dynamical Systems with Dynamic Stream Weights for Audiovisual Speaker Localization*, 2018

# Audiovisual speaker tracking

**Update step (Kalman filter with dynamic stream weights[1])**

Observation model:

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

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



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

[1] C. Schymura, T. Isenberg, D. Kolossa: *Extending Linear Dynamical Systems with Dynamic Stream Weights for Audiovisual Speaker Localization*, 2018

# Learning dynamic stream weights

**Standard approach: Supervised training with oracle dynamic stream weights**

Audio features → | Oracle DSW estimation | ← Transcription
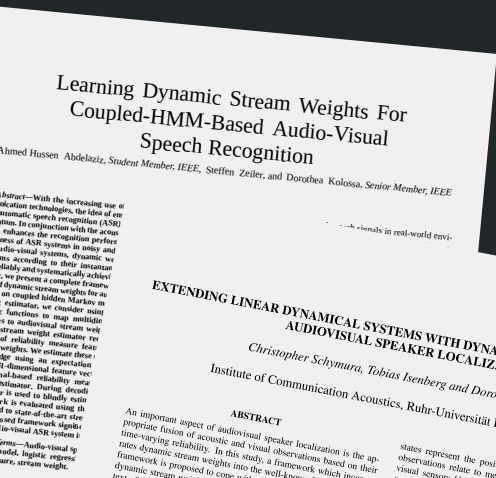Video features →

$\lambda^\star$

Reliability measures → | Parameter estimation $h(\mathbf{z}_k \mid \mathbf{w})$ | → $\mathbf{w}$

Learning Dynamic Stream Weights For Coupled-HMM-Based Audio-Visual Speech Recognition

Ahmed Hussen Abdelaziz, *Student Member, IEEE*, Steffen Zeiler, and Dorothea Kolossa, *Senior Member, IEEE*

EXTENDING LINEAR DYNAMICAL SYSTEMS WITH DYNAMIC STREAM WEIGHTS FOR AUDIOVISUAL SPEAKER LOCALIZATION
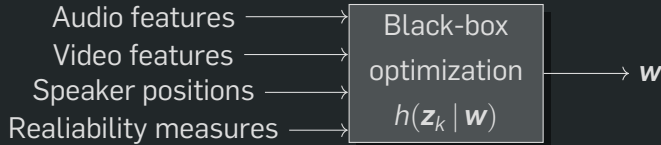
Christopher Schymura, Tobias Isenberg and Dorothea Kolossa

Institute of Communication Acoustics, Ruhr-Universität Bochum, Germany

# Learning dynamic stream weights

**Standard approach: Supervised training with oracle dynamic stream weights**

# Learning dynamic stream weights

**Proposed approach: Training with natural evolution strategies**

Audio features $\longrightarrow$
Video features $\longrightarrow$ Black-box
Speaker positions $\longrightarrow$ optimization
Realiability measures $\longrightarrow$ $h(\boldsymbol{z}_k \mid \boldsymbol{w})$ $\longrightarrow \boldsymbol{w}$

► No oracle information required.

► Flexible choice of loss/fitness function.

Daan Wierstra
Tom Schaul
DeepMind Technologies Ltd.
Fountain House, 130 Fenchurch Street
London, United Kingdom

Tobias Glasmachers
Institute for Neural Computation
Universitätsstrasse 150
Ruhr-University Bochum, Germany

Yi Sun
Google Inc.
1600 Amphitheatre Pkwy
Mountain View, United States

Jan Peters
Intelligent Autonomous Systems Institute
Hochschulstrasse 10
Technische Universität Darmstadt, Germany

Jürgen Schmidhuber
Dalle Molle di Studi sull'Intelligenza Artificiale (IDSIA)
...no (USI)/SUPSI

# 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
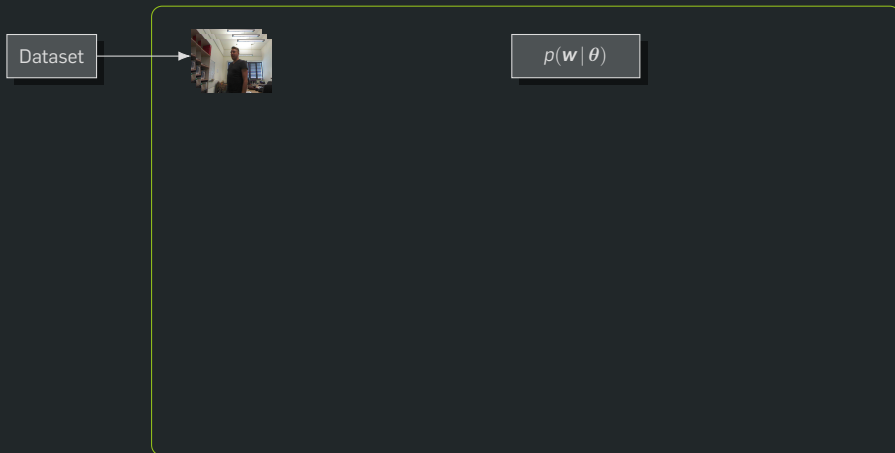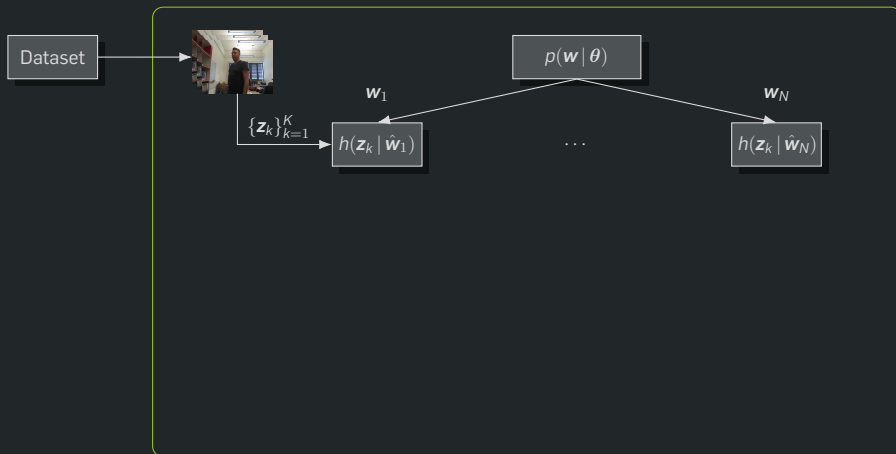
**Training procedure**

# Learning dynamic stream weights

**Training procedure**

# Learning dynamic stream weights

**Training procedure**

# Learning dynamic stream weights

**Implementation**

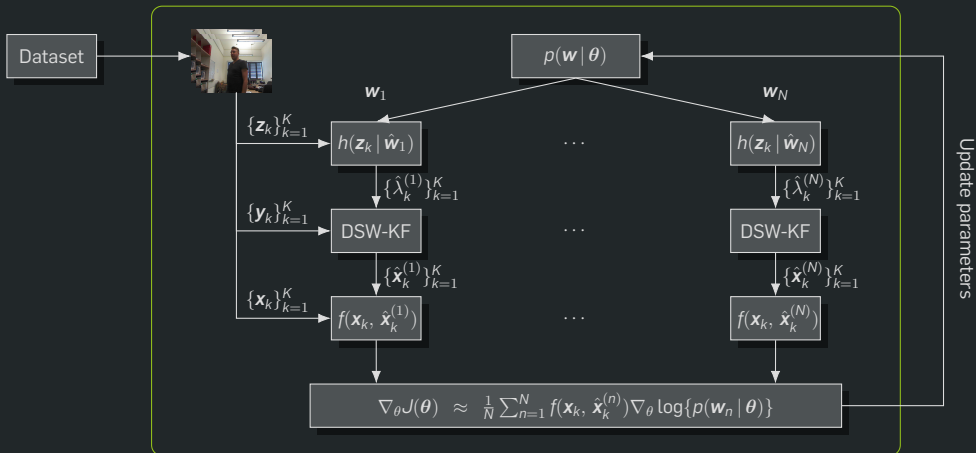- ▶ Reliability measures: instantaneous estimated a-priori SNR, acoustic and visual observation log-likelihoods[2].

[2] 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[2].

▶ Evaluation of two different DSW prediction models: logistic function and fully-connected feed-forward neural network.

---

[2] 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[2].

▶ Evaluation of two different DSW prediction models: logistic function and fully-connected feed-forward neural network.

▶ Separable natural evolution strategies (sNES) as optimizer:
$$p(\boldsymbol{w} \,|\, \boldsymbol{\theta}) = \mathcal{N}\Big(\boldsymbol{w} \,|\, \boldsymbol{\mu_w},\, \mathrm{diag}(\boldsymbol{\sigma_w})\Big)$$

[2] 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[2].

▶ Evaluation of two different DSW prediction models: logistic function and fully-connected feed-forward neural network.

▶ Separable natural evolution strategies (sNES) as optimizer:
$$p(\boldsymbol{w} \mid \boldsymbol{\theta}) = \mathcal{N}\Big(\boldsymbol{w} \mid \boldsymbol{\mu_w}, \, \mathrm{diag}(\boldsymbol{\sigma_w})\Big)$$

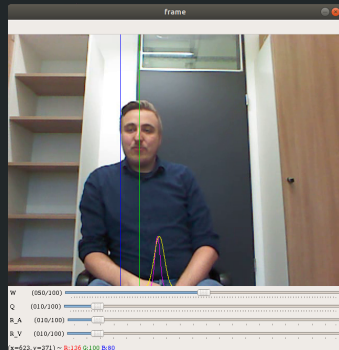▶ Fitness function allowing direct optimization of instantaneous localization error:
$$f(\boldsymbol{w}) = -\frac{1}{M} \sum_{m=1}^{M} \frac{1}{K_m} \sum_{k=1}^{K_m} \left( \phi_k^{(m)} - \hat{\phi}_k^{(m)}(\boldsymbol{w}) \right)^2$$

[2] 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[3] for acoustic localization, Viola-Jones[4] algorithm for visual localization.



---

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

[4] P. Viola, M. Jones: *Rapid object detection using a boosted cascade of simple features*, 2001

# Evaluation

## Experimental setup

▶ Front-end: DPD-MUSIC[3] for acoustic localization, Viola-Jones[4] algorithm for visual localization.

▶ Dataset of audiovisual recordings in an office environment ($T_{60} \approx 350\,\mathrm{ms}$) using the Kinect.
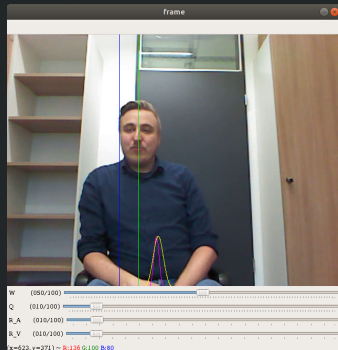
[3] Nadiri et al.: *Localization of multiple speakers under high reverberation using a spherical microphone array and the direct-path dominance test*, 2014

[4] P. Viola, M. Jones: *Rapid object detection using a boosted cascade of simple features*, 2001

# Evaluation

## Experimental setup

▶ Front-end: DPD-MUSIC[3] for acoustic localization, Viola-Jones[4] algorithm for visual localization.

▶ Dataset of audiovisual recordings in an office environment ($T_{60} \approx 350$ ms) using the Kinect.
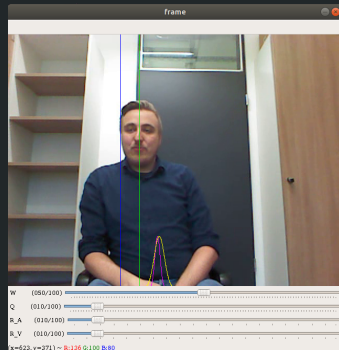
▶ Constant velocity dynamics model.



---

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

[4] P. Viola, M. Jones: *Rapid object detection using a boosted cascade of simple features*, 2001

# Evaluation

## Experimental setup

▶ Front-end: DPD-MUSIC[3] for acoustic localization, Viola-Jones[4] 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)
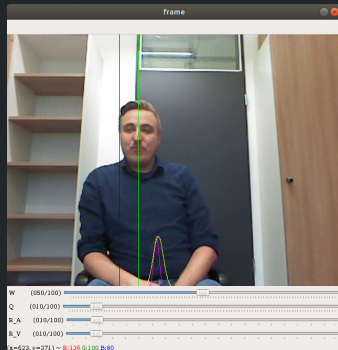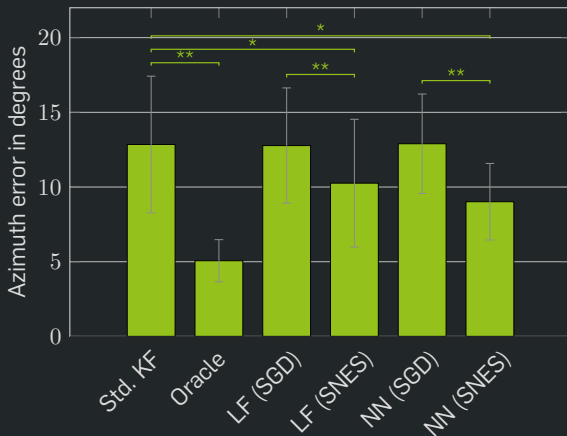


---

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

[4] P. Viola, M. Jones: *Rapid object detection using a boosted cascade of simple features*, 2001

# Evaluation

**Results**



Statistical significance: $\star$ for $p < 0.05$ and $\star\star$ 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.

**RUHR
UNIVERSITÄT
BOCHUM**

**RU**B

**Thank you for your attention!**