

Explainable Cognitive Models for Audiovisual Speaker Localization

12th AABBA General Meeting

Christopher Schymura and Dorothea Kolossa

16th January 2020

RUHR
UNIVERSITÄT
BOCHUM

RUB

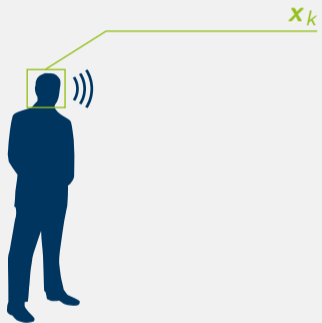
Problem statement



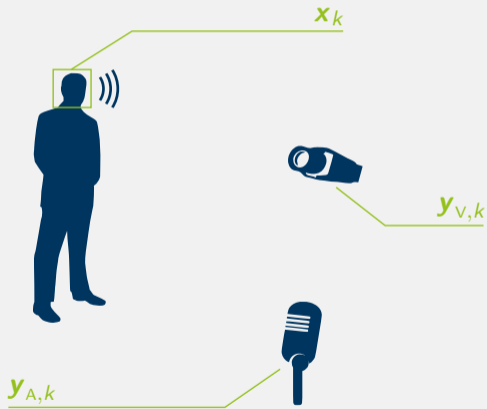
Problem statement



Problem statement



Problem statement

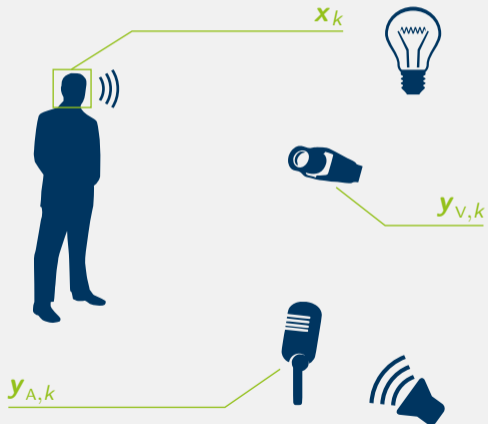


Observation functions:

$$y_{A,k} = h_A(x_k)$$

$$y_{V,k} = h_V(x_k)$$

Problem statement

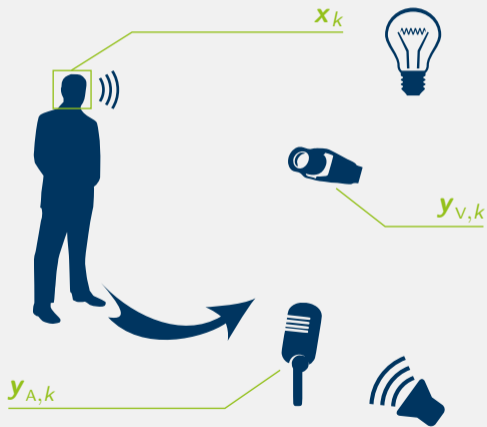


Observation functions:

$$\mathbf{y}_{A,k} = h_A(\mathbf{x}_k) + \mathbf{w}_{A,k}$$

$$\mathbf{y}_{V,k} = h_V(\mathbf{x}_k) + \mathbf{w}_{V,k}$$

Problem statement



State transition function:

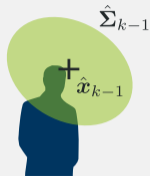
$$\mathbf{x}_k = f(\mathbf{x}_{k-1}) + \mathbf{v}_k$$

Observation functions:

$$\mathbf{y}_{A,k} = h_A(\mathbf{x}_k) + \mathbf{w}_{A,k}$$

$$\mathbf{y}_{V,k} = h_V(\mathbf{x}_k) + \mathbf{w}_{V,k}$$

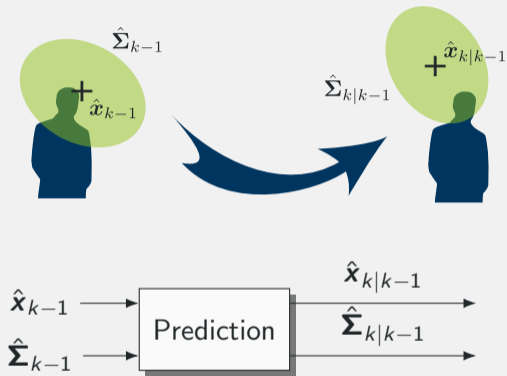
Recursive state estimation



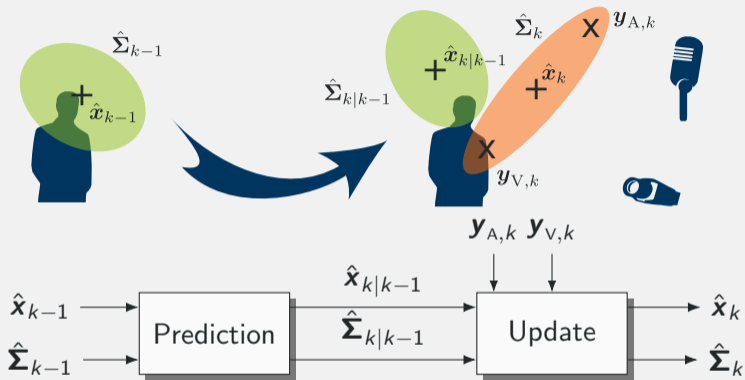
$$\hat{x}_{k-1}$$

$$\hat{\Sigma}_{k-1}$$

Recursive state estimation

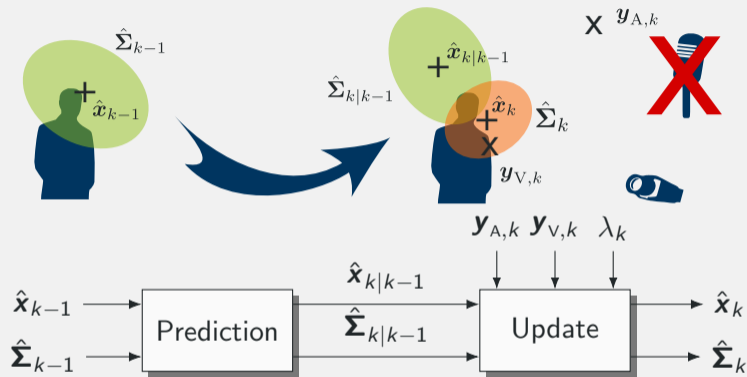


Recursive state estimation



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

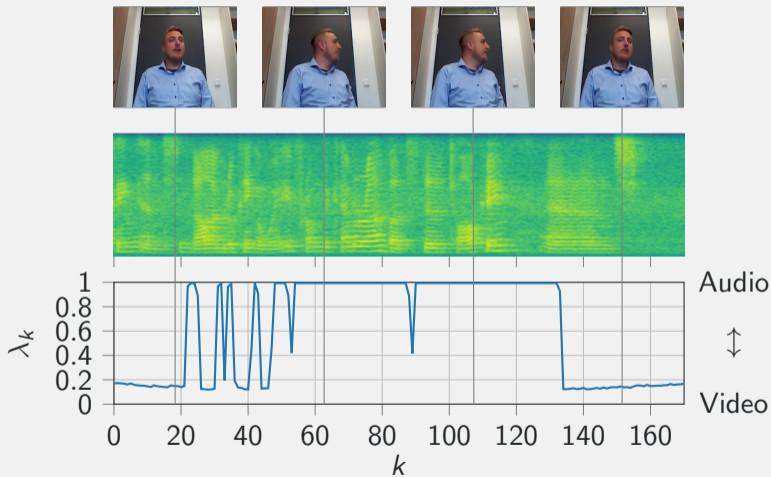
Recursive state estimation



$$\underbrace{p(\mathbf{x}_k | \mathbf{Y}_{A,1:k}, \mathbf{Y}_{V,1:k})}_{\text{Posterior}} \propto \underbrace{p(\mathbf{x}_k | \mathbf{Y}_{A,1:k-1}, \mathbf{Y}_{V,1:k-1})}_{\text{Prior}} \underbrace{p(\mathbf{y}_{A,k} | \mathbf{x}_k)^{\lambda_k} p(\mathbf{y}_{V,k} | \mathbf{x}_k)^{1-\lambda_k}}_{\text{Sensor model w. stream weights}^2}$$

²C. Schymura et al.: *Extending linear dynamical systems with dynamic stream weights for audiovisual speaker localization*, IWAENC, 2018

Dynamic stream weights



Dynamic stream weights

Assumption: $\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, k = 1, \dots, K$ fully observed, $\lambda_k \in [0, 1]$ and i.i.d.

Dynamic stream weights

Assumption: $\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, k = 1, \dots, K$ fully observed, $\lambda_k \in [0, 1]$ and i.i.d.

$$p(\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, \lambda_k) \propto p(\mathbf{y}_{A,k}|\mathbf{x}_k)^{\lambda_k} p(\mathbf{y}_{V,k}|\mathbf{x}_k)^{1-\lambda_k}$$
$$\Leftrightarrow \log\{p(\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, \lambda_k)\} = \lambda_k \log\{p(\mathbf{y}_{A,k}|\mathbf{x}_k)\} + (1 - \lambda_k) \log\{p(\mathbf{y}_{V,k}|\mathbf{x}_k)\} + c$$

Dynamic stream weights

Assumption: $\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, k = 1, \dots, K$ fully observed, $\lambda_k \in [0, 1]$ and i.i.d.

$$p(\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, \lambda_k) \propto p(\mathbf{y}_{A,k}|\mathbf{x}_k)^{\lambda_k} p(\mathbf{y}_{V,k}|\mathbf{x}_k)^{1-\lambda_k}$$
$$\Leftrightarrow \log\{p(\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, \lambda_k)\} = \lambda_k \log\{p(\mathbf{y}_{A,k}|\mathbf{x}_k)\} + (1 - \lambda_k) \log\{p(\mathbf{y}_{V,k}|\mathbf{x}_k)\} + c$$

Problem: Direct optimization not feasible.

Dynamic stream weights

Assumption: $\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, k = 1, \dots, K$ fully observed, $\lambda_k \in [0, 1]$ and i.i.d.

$$p(\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, \lambda_k) \propto p(\mathbf{y}_{A,k}|\mathbf{x}_k)^{\lambda_k} p(\mathbf{y}_{V,k}|\mathbf{x}_k)^{1-\lambda_k}$$
$$\Leftrightarrow \log\{p(\mathbf{x}_k, \mathbf{y}_{A,k}, \mathbf{y}_{V,k}, \lambda_k)\} = \lambda_k \log\{p(\mathbf{y}_{A,k}|\mathbf{x}_k)\} + (1 - \lambda_k) \log\{p(\mathbf{y}_{V,k}|\mathbf{x}_k)\} + c$$

Problem: Direct optimization not feasible.

Solution: Impose prior on λ_k , e.g. Gaussian or symmetric Beta³ distribution.

$$J(\lambda_k) = \lambda_k \log\{p(\mathbf{y}_{A,k}|\mathbf{x}_k)\} + (1 - \lambda_k) \log\{p(\mathbf{y}_{V,k}|\mathbf{x}_k)\} + \log\{p(\lambda_k)\}$$

³C. Schymura et al.: *Audiovisual speaker tracking using nonlinear dynamical systems with dynamic stream weights*, arXiv, 2019

Symmetric Beta prior

$$J(\lambda_k) = \lambda_k \log\{p(\mathbf{y}_{A,k}|\mathbf{x}_k)\} + (1 - \lambda_k) \log\{p(\mathbf{y}_{V,k}|\mathbf{x}_k)\} + \log\{p(\lambda_k)\}$$

with

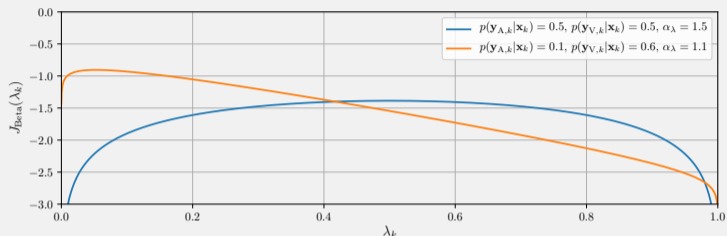
$$p(\lambda_k) = \frac{1}{B(\alpha_\lambda, \alpha_\lambda)} \lambda_k^{\alpha_\lambda-1} (1 - \lambda_k)^{\alpha_\lambda-1}$$

yields

$$\begin{aligned} J_{\text{Beta}}(\lambda_k) &= \lambda_k \log\{p(\mathbf{y}_{A,k}|\mathbf{x}_k)\} + (1 - \lambda_k) \log\{p(\mathbf{y}_{V,k}|\mathbf{x}_k)\} \\ &\quad + (\alpha_\lambda - 1) \left(\log\{\lambda_k\} + \log\{1 - \lambda_k\} \right) + \text{const.} \end{aligned}$$

$$\Rightarrow \lambda_k^* = \max_{\lambda_k} J_{\text{Beta}}(\lambda_k) \quad \text{s. t.} \quad 0 < \lambda_k < 1$$

Symmetric Beta prior

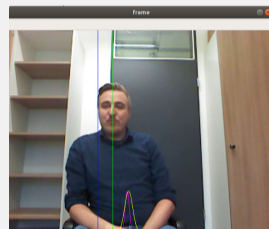
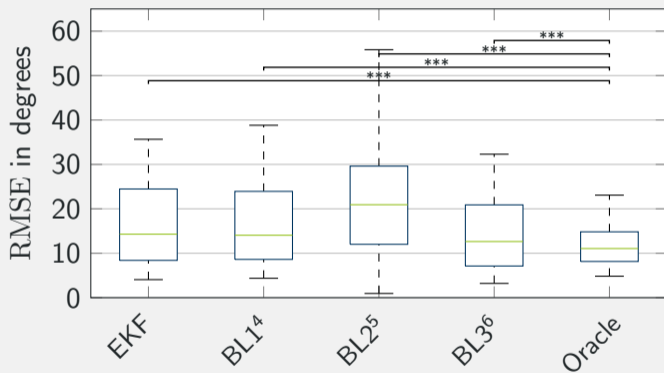


$J_{\text{Beta}}(\lambda_k)$ is a concave function:

$$\frac{dJ_{\text{Beta}}(\lambda_k)}{d\lambda_k} = \log \left\{ \frac{p(\mathbf{y}_{A,k}|\mathbf{x}_k)}{p(\mathbf{y}_{V,k}|\mathbf{x}_k)} \right\} + (\alpha_\lambda - 1) \left(\frac{1}{\lambda_k} + \frac{1}{\lambda_k - 1} \right)$$

$$\frac{d^2 J_{\text{Beta}}(\lambda_k)}{d\lambda_k^2} = (\alpha_\lambda - 1) \left(\frac{1}{\lambda_k^2} + \frac{1}{(\lambda_k - 1)^2} \right) < 0 \quad \forall \alpha_\lambda > 1$$

Results I



[* * *] $p < 0.001$

⁴T. Gehrig et al.: *Kalman filters for audio-video source localization*, WASPAA, 2005

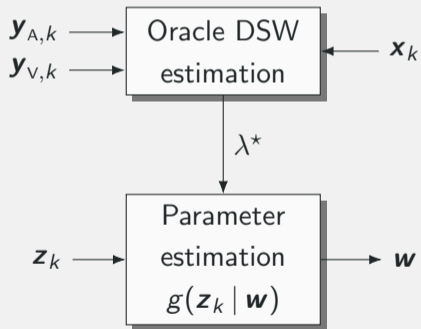
⁵S. Gerlach et al.: *2D audio-visual localization in home environments using a particle filter*, ITG Symp., 2012

⁶X. Qian et al.: *3D audio-visual speaker tracking with an adaptive particle filter*, ICASSP, 2017

Learning dynamic stream weights

Supervised learning approach

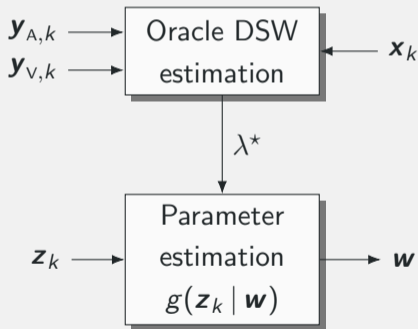
Oracle DSW serve as targets



Learning dynamic stream weights

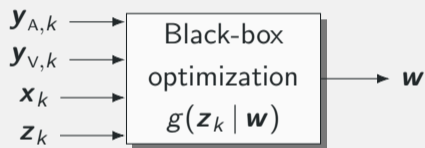
Supervised learning approach

Oracle DSW serve as targets



Evolutionary⁷ learning approach

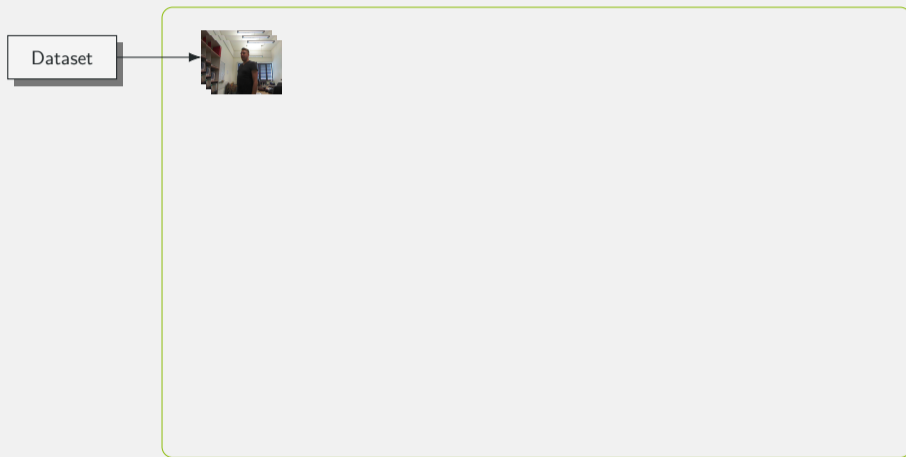
Direct optimization of localization error



⁷D. Wierstra et al.: *Natural evolution strategies*, Journal of machine learning research, vol. 15, 2014

Learning dynamic stream weights

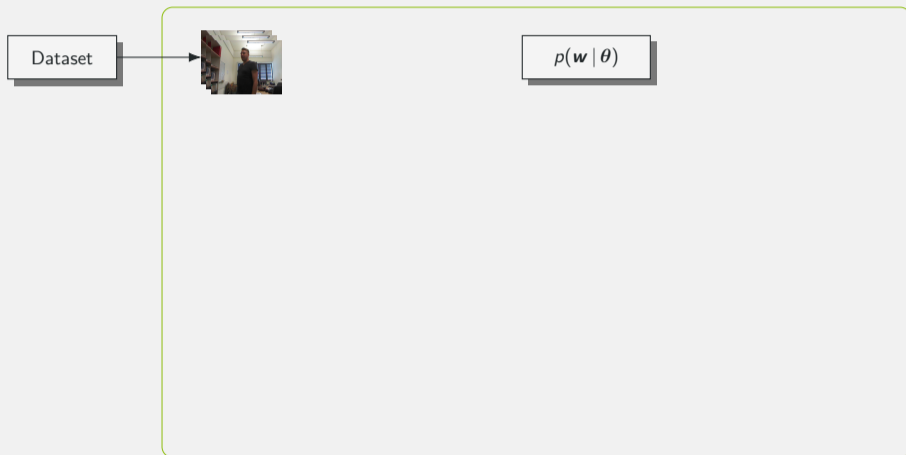
Training procedure⁸



⁸C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

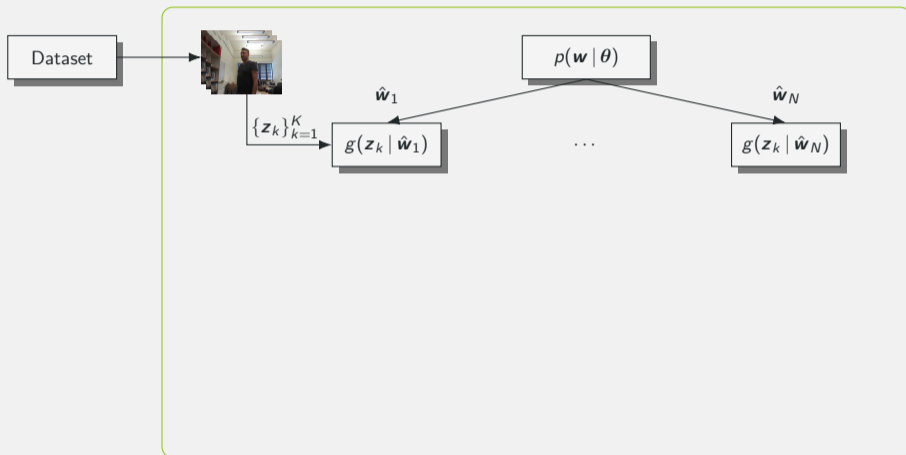
Training procedure⁸



⁸C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

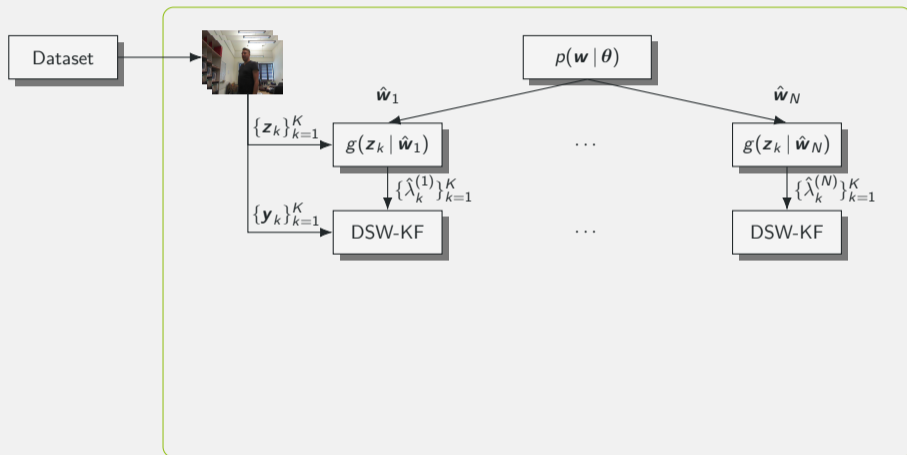
Training procedure⁸



⁸C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

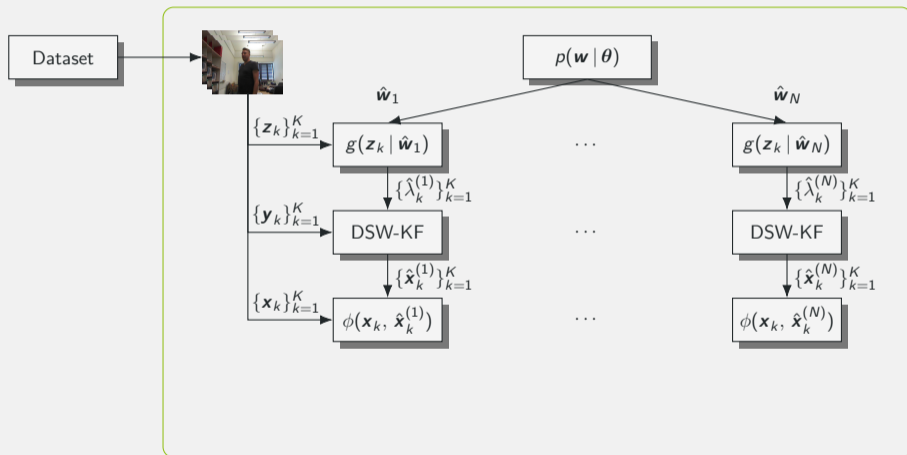
Training procedure⁸



⁸C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

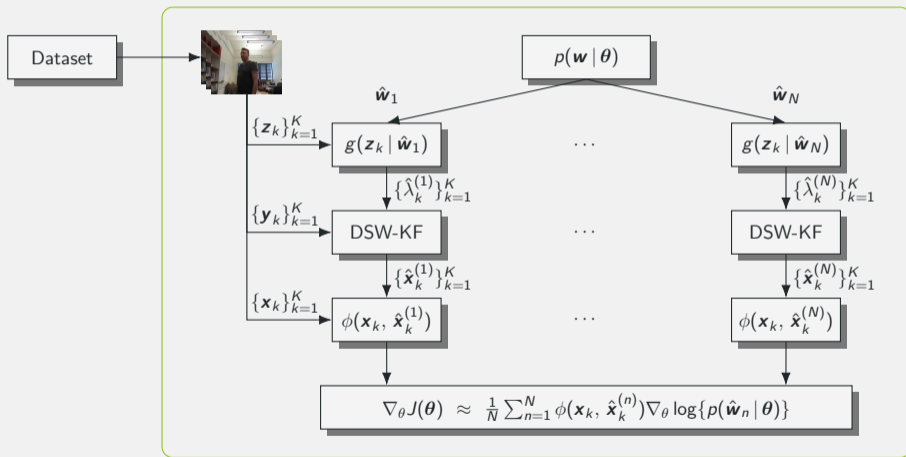
Training procedure⁸



⁸C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Learning dynamic stream weights

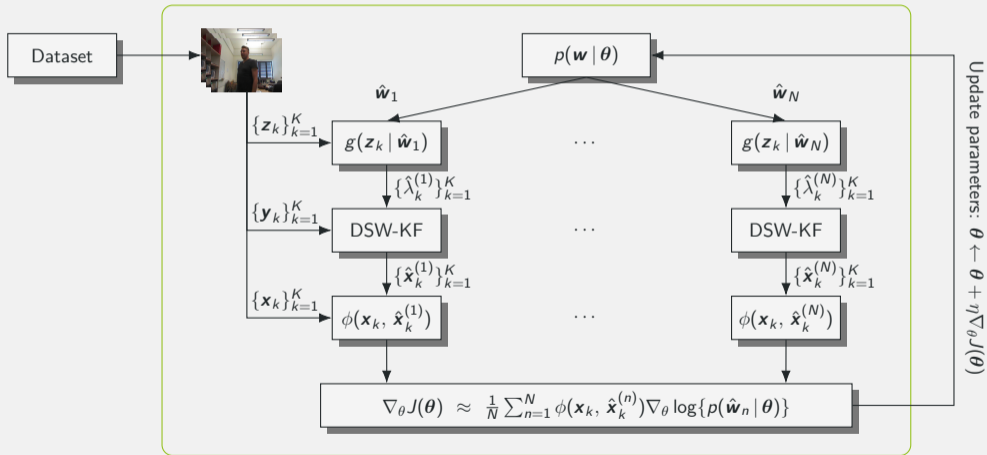
Training procedure⁸



⁸C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

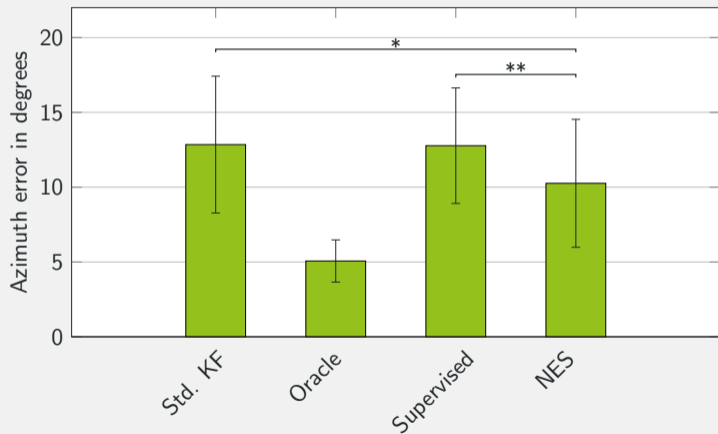
Learning dynamic stream weights

Training procedure⁸



⁸C. Schymura et al.: *Learning dynamic stream weights for linear dynamical systems using natural evolution strategies*, ICASSP, 2019

Results II

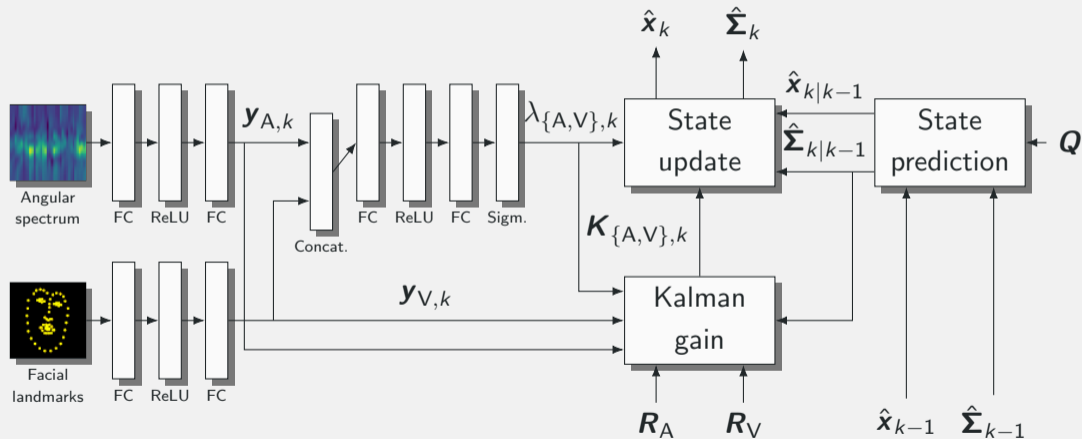


[*] $p < 0.05$

[**] $p < 0.01$

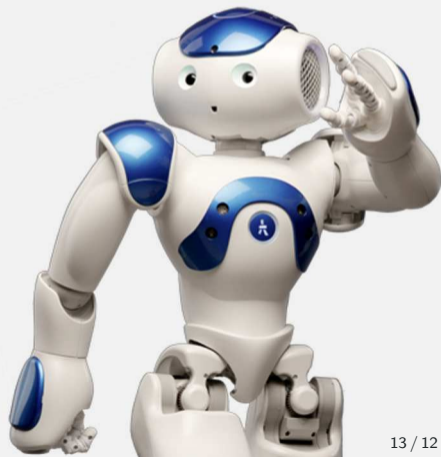
Ongoing work (in collaboration with NTT)

End-to-end optimization in a deep learning framework:



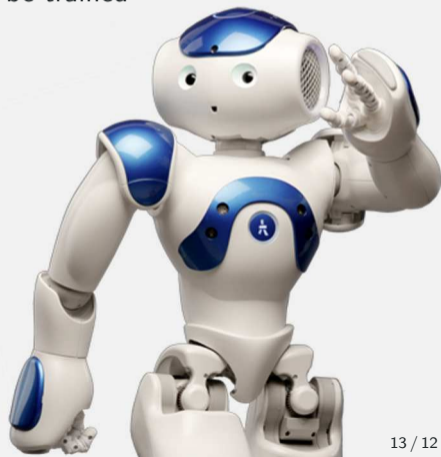
Conclusions

- ▶ Dynamic stream weights can benefit audiovisual speaker localization performance and provide an **additional notion of uncertainty**.



Conclusions

- ▶ Dynamic stream weights can benefit audiovisual speaker localization performance and provide an **additional notion of uncertainty**.
- ▶ Models for estimating dynamic stream weights can be trained effectively using **evolutionary optimization**.



Conclusions

- ▶ Dynamic stream weights can benefit audiovisual speaker localization performance and provide an **additional notion of uncertainty**.
- ▶ Models for estimating dynamic stream weights can be trained effectively using **evolutionary optimization**.
- ▶ State error covariance matrices and dynamic stream weights jointly contribute to **model explainability**.



Conclusions

- ▶ Dynamic stream weights can benefit audiovisual speaker localization performance and provide an **additional notion of uncertainty**.
- ▶ Models for estimating dynamic stream weights can be trained effectively using **evolutionary optimization**.
- ▶ State error covariance matrices and dynamic stream weights jointly contribute to **model explainability**.
- ▶ Promising direction for future research: How to integrate modern deep learning techniques without sacrificing explainability?



Conclusions

- ▶ Dynamic stream weights can benefit audiovisual speaker localization performance and provide an **additional notion of uncertainty**.
- ▶ Models for estimating dynamic stream weights can be trained effectively using **evolutionary optimization**.
- ▶ State error covariance matrices and dynamic stream weights jointly contribute to **model explainability**.
- ▶ Promising direction for future research: How to integrate modern deep learning techniques without sacrificing explainability?

Thank you for your attention!

